

Classification (overview)

Probability and Statistics for Data Science

Carlos Fernandez-Granda

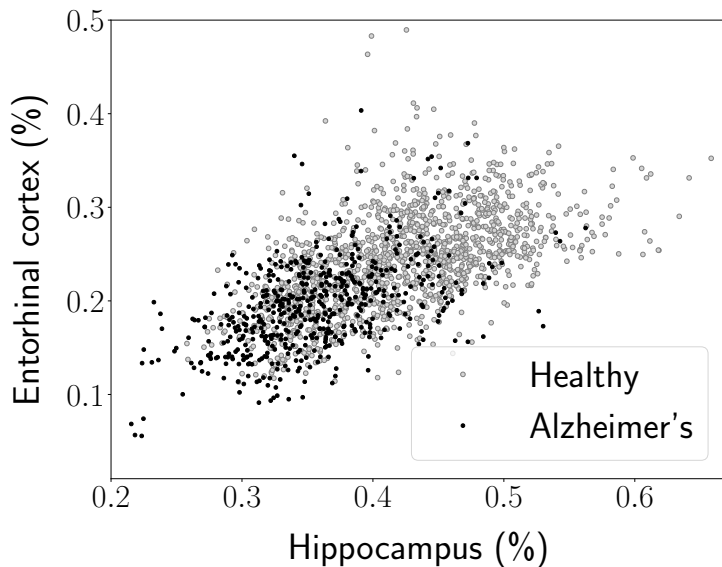


These slides are based on the book [Probability and Statistics for Data Science](#) by Carlos Fernandez-Granda, available for purchase [here](#). A free preprint, videos, code, slides and solutions to exercises are available at <https://www.ps4ds.net>

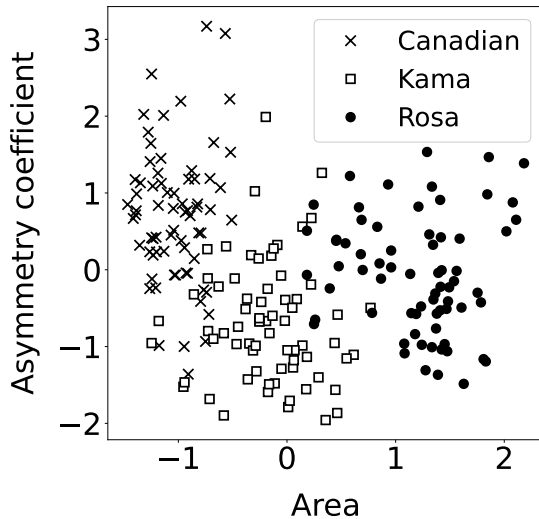
Classification

Goal: Assign class to data based on features

Alzheimer's diagnostics



Identification of wheat varieties



U.S. House of Representatives in 1984

Predict political affiliation (Republican or Democrat) from voting record on 16 issues

Probabilistic modeling

Class: Random variable \tilde{y} with range Y

Features: Random vector \tilde{x} with range \mathcal{X}

Maximum a posteriori (MAP) estimator of \tilde{y} given $\tilde{x} = x \in \mathcal{X}$

$$\text{MAP}(x) := \arg \max_{y \in Y} p_{\tilde{y} | \tilde{x}}(y | x)$$

MAP estimation is optimal

For any other estimator h

$$\begin{aligned} P(\tilde{y} = h(\tilde{x})) &= \sum_{x \in \mathcal{X}} P(\tilde{x} = x, \tilde{y} = h(\tilde{x})) \\ &= \sum_{x \in \mathcal{X}} P(\tilde{x} = x) P(\tilde{y} = h(x) \mid \tilde{x} = x) \\ &= \sum_{x \in \mathcal{X}} p_{\tilde{x}}(x) p_{\tilde{y} \mid \tilde{x}}(h(x) \mid x) \\ &\leq \sum_{x \in \mathcal{X}} p_{\tilde{x}}(x) p_{\tilde{y} \mid \tilde{x}}(\text{MAP}(x) \mid x) \\ &= \sum_{x \in \mathcal{X}} P(\tilde{x} = x) P(\tilde{y} = \text{MAP}(x) \mid \tilde{x} = x) \\ &= P(\tilde{y} = \text{MAP}(\tilde{x})) \end{aligned}$$

Are we done here?

Prediction of political affiliation

16 binary features (Yes/No votes)

Can we estimate $p_{\tilde{y}|\tilde{x}}(\cdot|x)$ for any x ? **No**

Possible values of x ? $2^{16} = 65,536$

We have 425 training examples...

We need assumptions to build tractable approximations

Plan

Generative Models

Discriminative Models

Evaluation

Generative Models

Discriminative Models

Evaluation

Generative classification models

Goal: Estimate conditional pmf $p_{\tilde{y}|\tilde{x}}$ of class given features

1. Use data to approximate

- ▶ Pmf $p_{\tilde{y}}$ of class
- ▶ Conditional pmf $p_{\tilde{x}|\tilde{y}}$ or conditional pdf $f_{\tilde{x}|\tilde{y}}$ of features given class

2. Apply Bayes' rule

$$p_{\tilde{y}|\tilde{x}}(y|x) = \frac{p_{\tilde{y}}(y)p_{\tilde{x}|\tilde{y}}(x|y)}{p_{\tilde{x}}(x)} \quad \text{or} \quad \frac{p_{\tilde{y}}(y)f_{\tilde{x}|\tilde{y}}(x|y)}{f_{\tilde{x}}(x)}$$

Challenge

Approximating conditional (joint) pmf / pdf of features given class

- ▶ Naive Bayes:

Features are **conditionally independent** given class

- ▶ Gaussian discriminant analysis:

Features are conditionally **Gaussian** given class

Prediction of political affiliation

Class

$$\tilde{y} = \begin{cases} R & \text{Republican} \\ D & \text{Democrat} \end{cases}$$

Features ($1 \leq i \leq 16$)

$$\tilde{x}[i] = \begin{cases} 1 & \text{voted Yes on issue } i \\ 0 & \text{otherwise} \end{cases}$$

Naive Bayes

We assume votes are conditionally independent given affiliation

$$p_{\tilde{x}|\tilde{y}}(x|R) = \prod_{i=1}^d p_{\tilde{x}[i]|\tilde{y}}(x[i]|R)$$
$$p_{\tilde{x}|\tilde{y}}(x|D) = \prod_{i=1}^d p_{\tilde{x}[i]|\tilde{y}}(x[i]|D)$$

Bayes rule

$$\begin{aligned} & p_{\tilde{y}|\tilde{x}}(R|x) \\ &= \frac{p_{\tilde{y},\tilde{x}}(R,x)}{p_{\tilde{x}}(x)} \\ &= \frac{p_{\tilde{y}}(R)p_{\tilde{x}|\tilde{y}}(x|R)}{p_{\tilde{y},\tilde{x}}(R,x) + p_{\tilde{y},\tilde{x}}(D,x)} \\ &= \frac{p_{\tilde{y}}(R) \prod_{i=1}^d p_{\tilde{x}[i]|\tilde{y}}(x[i]|R)}{p_{\tilde{y}}(R) \prod_{i=1}^d p_{\tilde{x}[i]|\tilde{y}}(x[i]|R) + p_{\tilde{y}}(D) \prod_{i=1}^d p_{\tilde{x}[i]|\tilde{y}}(x[i]|D)} \end{aligned}$$

Estimated probabilities

$$p_{\tilde{y}}(R) = 0.381 \quad (p_{\tilde{y}}(D) = 0.619)$$

i	1	2	3	4	5	6	7	8
$p_{\tilde{x}[i] \tilde{y}}(1 R)$	0.19	0.50	0.14	0.99	0.95	0.90	0.24	0.15
$p_{\tilde{x}[i] \tilde{y}}(1 D)$	0.61	0.50	0.89	0.05	0.22	0.47	0.78	0.83

i	9	10	11	12	13	14	15	16
$p_{\tilde{x}[i] \tilde{y}}(1 R)$	0.11	0.55	0.14	0.87	0.86	0.98	0.09	0.66
$p_{\tilde{x}[i] \tilde{y}}(1 D)$	0.76	0.47	0.51	0.15	0.29	0.35	0.64	0.94

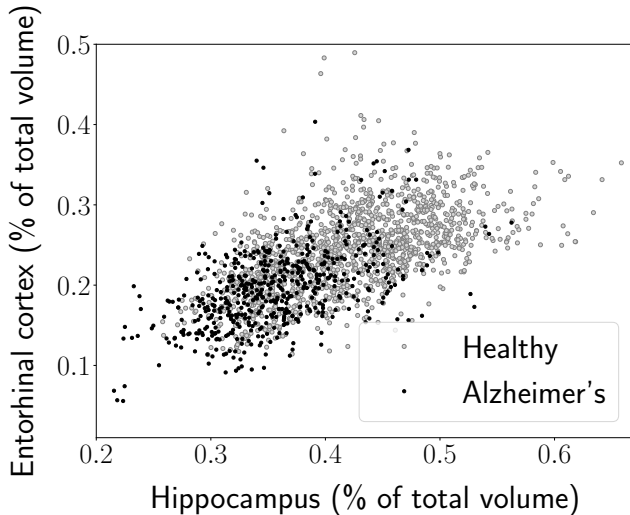
Applying the model

i	1	2	3	4	5	6	7	8
$p_{\tilde{x}[i] \tilde{y}}(1 R)$	0.19	0.50	0.14	0.99	0.95	0.90	0.24	0.15
$p_{\tilde{x}[i] \tilde{y}}(1 D)$	0.61	0.50	0.89	0.05	0.22	0.47	0.78	0.83
Example	N	–	Y	N	N	Y	Y	Y

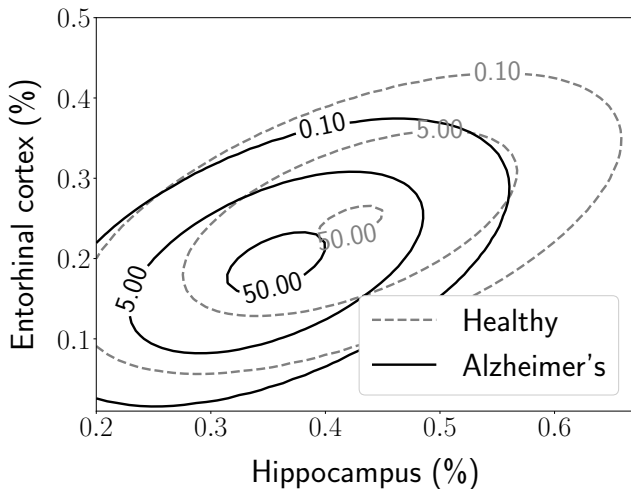
i	9	10	11	12	13	14	15	16
$p_{\tilde{x}[i] \tilde{y}}(1 R)$	0.11	0.55	0.14	0.87	0.86	0.98	0.09	0.66
$p_{\tilde{x}[i] \tilde{y}}(1 D)$	0.76	0.47	0.51	0.15	0.29	0.35	0.64	0.94
Example	N	Y	N	N	N	N	Y	–

$$\begin{aligned}
 & p_{\tilde{y} | \tilde{x}}(D | x) \\
 &= \frac{p_{\tilde{y}}(D) \prod_{i \in \{1,3,\dots,15\}} p_{\tilde{x}[i] | \tilde{y}}(x[i] | D)}{p_{\tilde{y}}(D) \prod_{i \in \{1,3,\dots,15\}} p_{\tilde{x}[i] | \tilde{y}}(x[i] | D) + p_{\tilde{y}}(R) \prod_{i \in \{1,3,\dots,15\}} p_{\tilde{x}[i] | \tilde{y}}(x[i] | R)} \\
 &= 1 - 1.410^{-8}
 \end{aligned}$$

Gaussian discriminant analysis



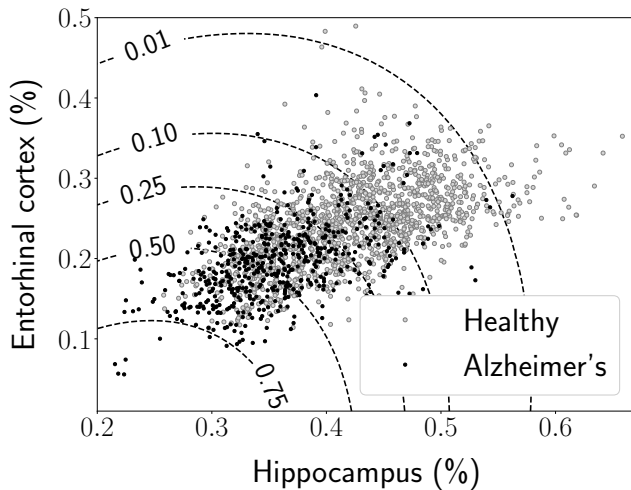
Gaussian parametric model of features given class



Bayes' rule

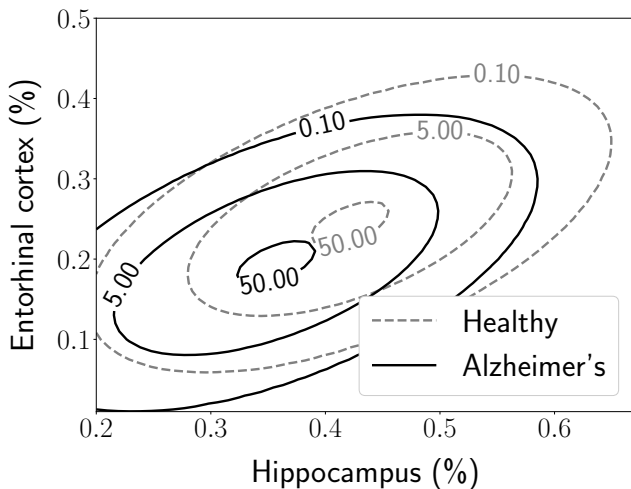
$$\begin{aligned} p_{\tilde{y} | \tilde{x}}(y | x) &= \frac{p_{\tilde{y}}(y) f_{\tilde{x} | \tilde{y}}(x | y)}{\sum_{k \in \{1, 2, \dots, c\}} p_{\tilde{y}}(k) f_{\tilde{x} | \tilde{y}}(x | k)} \\ &= \frac{\frac{p_{\tilde{y}}(y)}{\sqrt{(2\pi)^d |\Sigma_y|}} \exp\left(-\frac{1}{2} (x - \mu_y)^T \Sigma_y^{-1} (x - \mu_y)\right)}{\sum_{k \in \{1, 2, \dots, c\}} \frac{p_{\tilde{y}}(k)}{\sqrt{(2\pi)^d |\Sigma_k|}} \exp\left(-\frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k)\right)} \end{aligned}$$

$$p_{\tilde{y}|\tilde{x}}$$

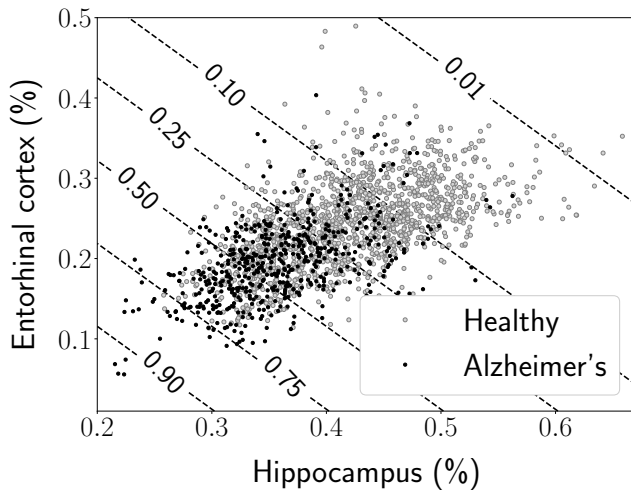


Linear discriminant analysis

We fit $\Sigma_a = \Sigma_b = \Sigma$



Linear discriminant analysis



Generative Models

Discriminative Models

Evaluation

Discriminative models

Goal: Directly approximate conditional pmf $p_{\tilde{y}|\tilde{x}}$ of class given features

Different assumptions lead to different models p_{Θ}

- ▶ **Linear:** Logistic / softmax regression
- ▶ **Nonlinear:** Neural networks / classification trees

Parameters Θ are learned from data

Data: $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$

Model the i th feature and label as the random variables \tilde{x}_i and \tilde{y}_i

Maximize the **conditional** likelihood of the labels given the features

Conditional likelihood and log-likelihood

Assumption 1:

Labels are conditionally independent given the features

Assumption 2:

\tilde{y}_i is conditionally independent from $\{\tilde{x}_m\}_{m \neq i}$ given \tilde{x}_i

$$\begin{aligned}\mathcal{L}_{XY}(\Theta) &:= P(\tilde{y}_1 = y_1, \dots, \tilde{y}_n = y_n \mid \tilde{x}_1 = x_1, \dots, \tilde{x}_n = x_n) \\ &= \prod_{k=1}^c \prod_{\{i: y_i=k\}} p_{\Theta}(x_i)_k \\ \log \mathcal{L}_{XY}(\Theta) &= \sum_{k=1}^c \sum_{\{i: y_i=k\}} \log p_{\Theta}(x_i)_k\end{aligned}$$

Can be difficult to maximize!

Generalized linear models

Idea: Approximate $p_{\tilde{y}|\tilde{x}}$ as a linear function of the features

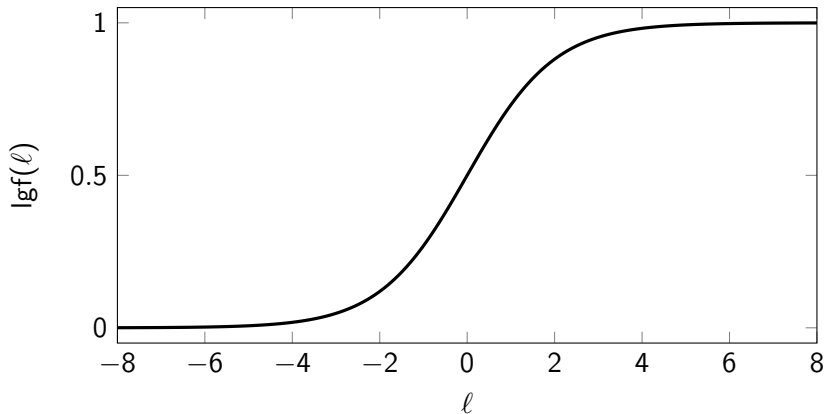
Problem: Linear functions are *not valid probabilities*

Solution: Map linear function to $[0,1]$ using **link function**:

- ▶ Logistic function for binary classification
- ▶ Softmax function for multiclass classification

Logistic function

$$\text{lgf}(\ell) := \frac{\exp(\ell)}{1 + \exp(\ell)} = \frac{1}{1 + \exp(-\ell)}$$

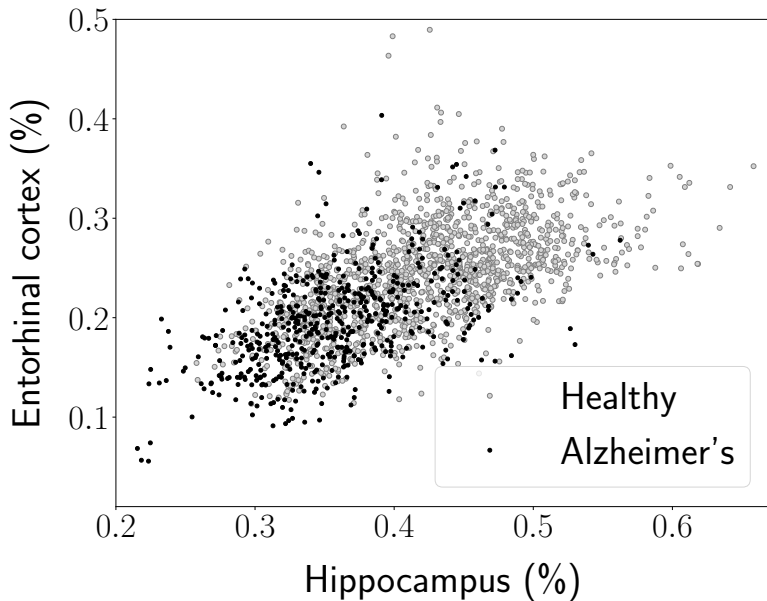


Logistic regression

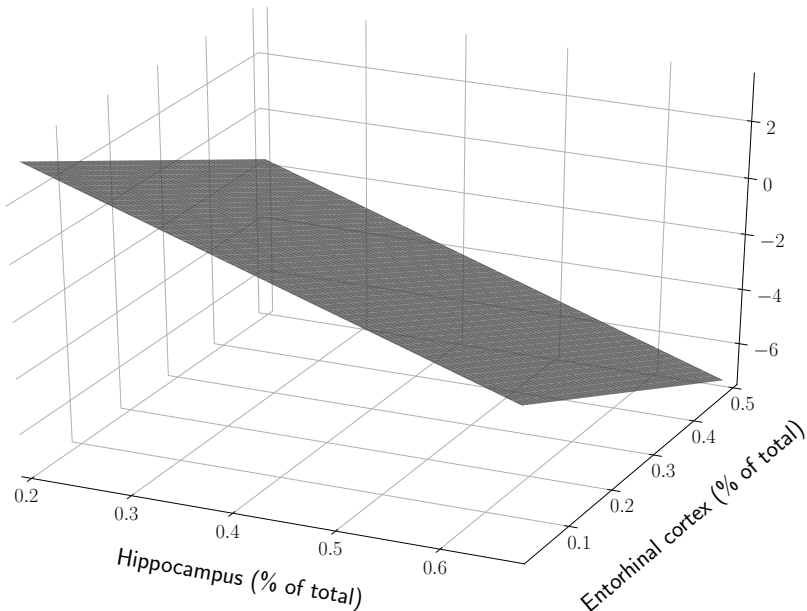
$$P(\tilde{y} = 1 \mid \tilde{x} = x) \approx p_{\Theta}(x) := \text{lgf}(\beta^T x + \alpha)$$

Concave log-likelihood maximized via iterative methods

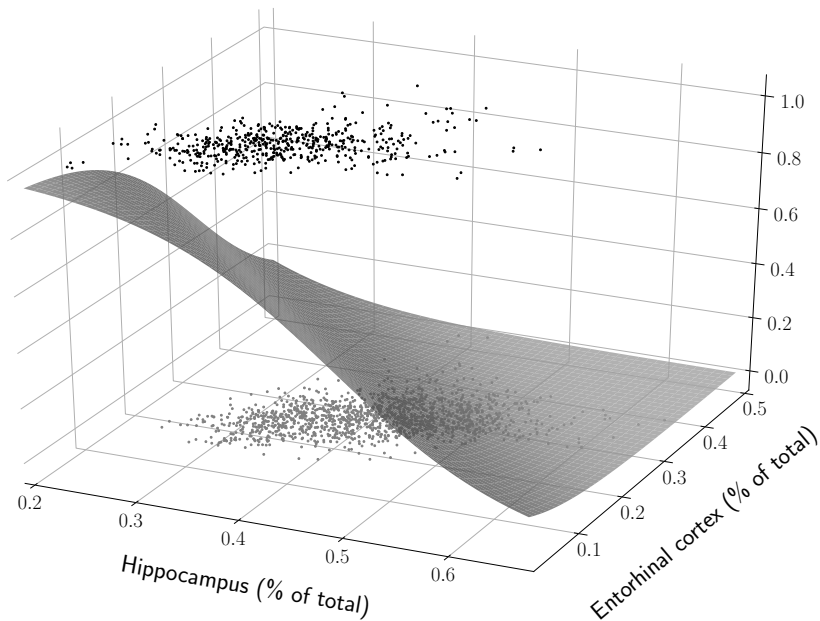
Alzheimer's diagnosis



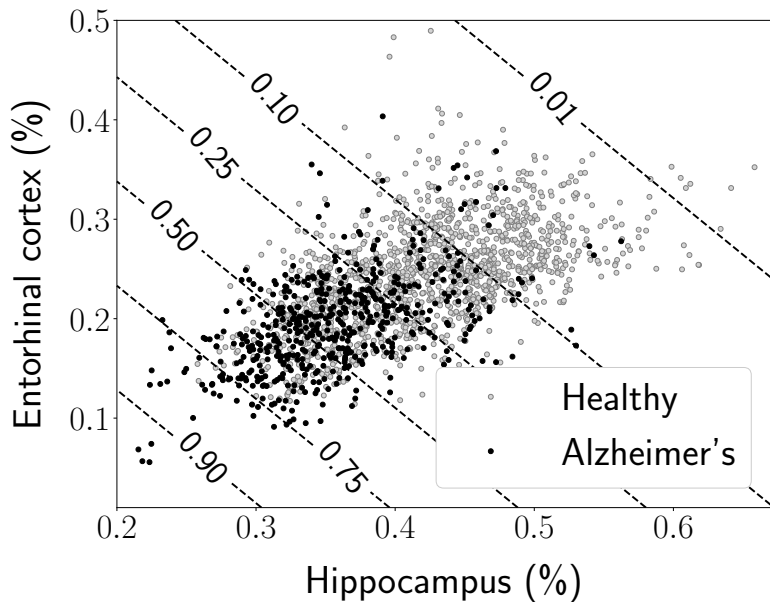
$$-11.9 x_{\text{hippocampus}} - 10.5 x_{\text{entorhinal}} + 5.9$$



$$\lgf(-11.9 x_{\text{hippocampus}} - 10.5 x_{\text{entorhinal}} + 5.9)$$



$$\lgf(-11.9 x_{\text{hippocampus}} - 10.5 x_{\text{entorhinal}} + 5.9)$$

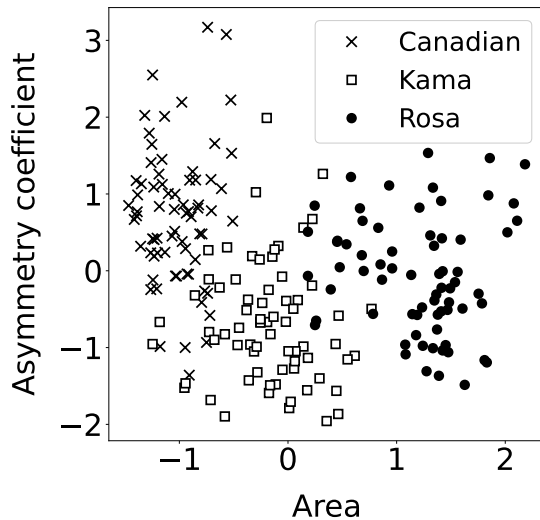


Softmax

$$P(\tilde{y} = k \mid \tilde{x} = x) \approx p_{\Theta}(x) := \frac{\exp(\beta_k^T x + \alpha_k)}{\sum_{l=1}^c \exp(\beta_l^T x + \alpha_l)} \quad 1 \leq k \leq c$$

Concave log-likelihood maximized via iterative methods

Wheat varieties



Softmax acts like a soft maximum

$$x_{\text{area}} := -1, x_{\text{asym}} := 2$$

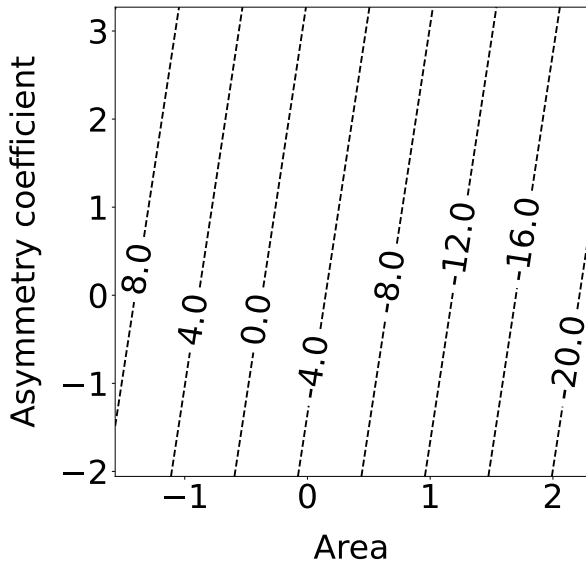
$$\begin{bmatrix} \text{Canadian} \\ \text{Kama} \\ \text{Rosa} \end{bmatrix} = \begin{bmatrix} -7.7 x_{\text{area}} + 0.9 x_{\text{asym}} - 2.9 \\ 0.4 x_{\text{area}} - 1.2 x_{\text{asym}} + 2.7 \\ 7.3 x_{\text{area}} + 0.4 x_{\text{asym}} + 0.2 \end{bmatrix} = \begin{bmatrix} 6.6 \\ -0.1 \\ -6.3 \end{bmatrix}$$

$$\exp(6.6) = 735 \quad \exp(-0.1) = 0.905 \quad \exp(-6.3) = 0.002$$

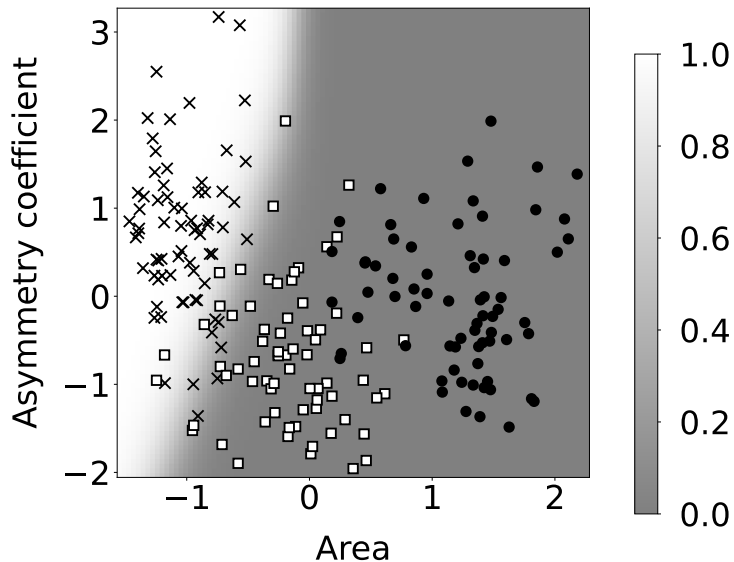
Softmax:

$$\begin{bmatrix} P(\text{Canadian} \mid x_{\text{area}}, x_{\text{asym}}) \\ P(\text{Kama} \mid x_{\text{area}}, x_{\text{asym}}) \\ P(\text{Rosa} \mid x_{\text{area}}, x_{\text{asym}}) \end{bmatrix} = \begin{bmatrix} \frac{735}{735+0.905+0.002} \\ \frac{0.905}{735+0.905+0.002} \\ \frac{0.002}{735+0.905+0.002} \end{bmatrix} = \begin{bmatrix} 0.999 \\ 0.001 \\ 0.000 \end{bmatrix}$$

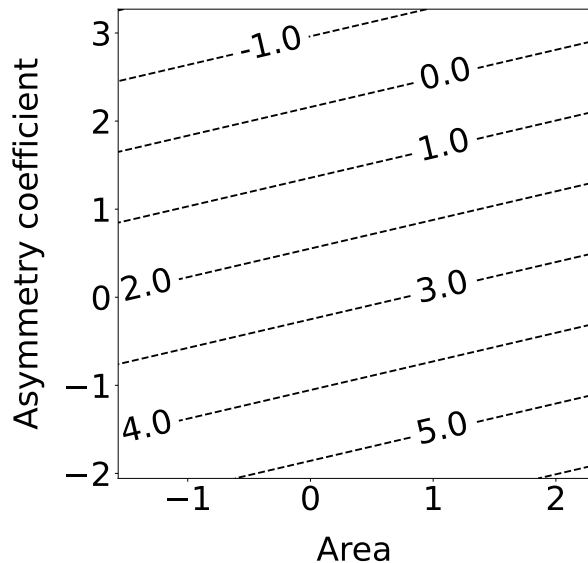
Canadian: $-7.7 x_{\text{area}} + 0.9 x_{\text{asym}} - 2.9$



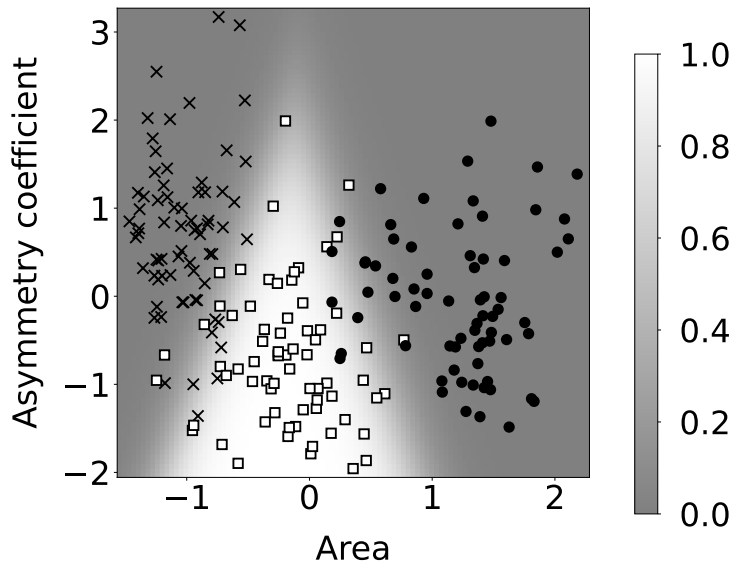
Canadian: Estimated probability



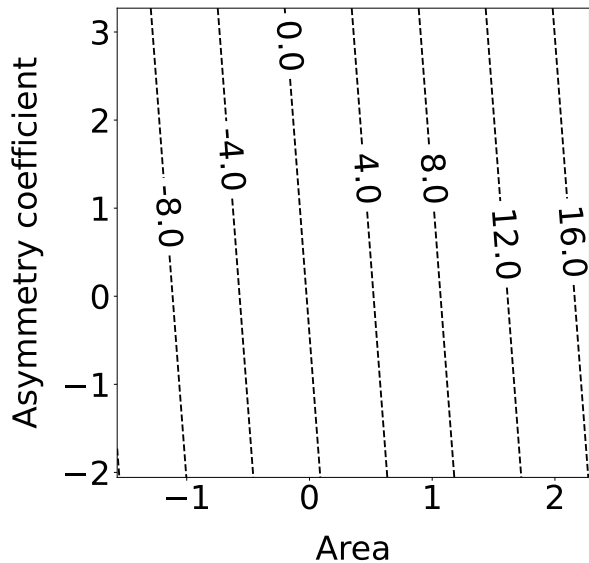
Kama: $0.4 x_{\text{area}} - 1.2 x_{\text{asym}} + 2.7$



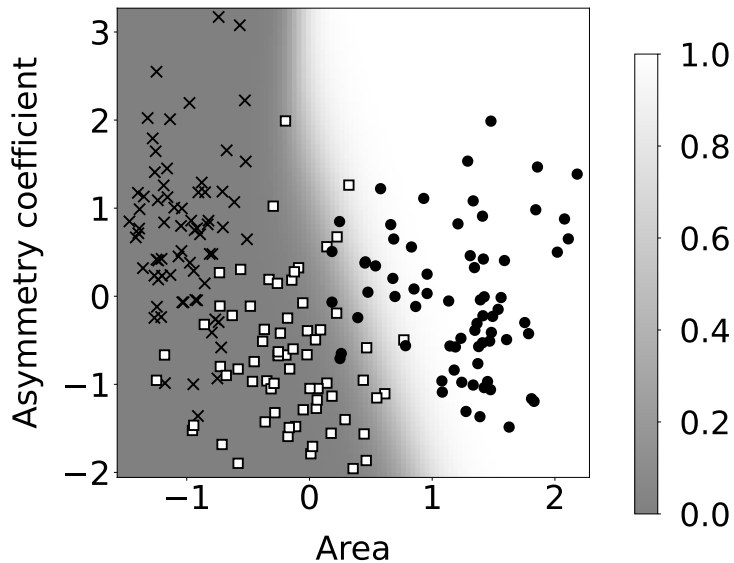
Kama: Estimated probability



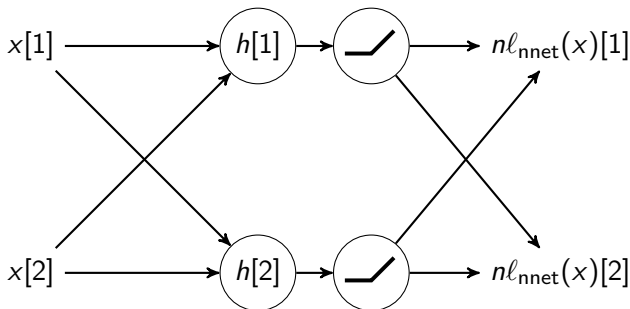
Rosa: $7.3 x_{\text{area}} + 0.4 x_{\text{asym}} + 0.2$



Rosa: Estimated probability



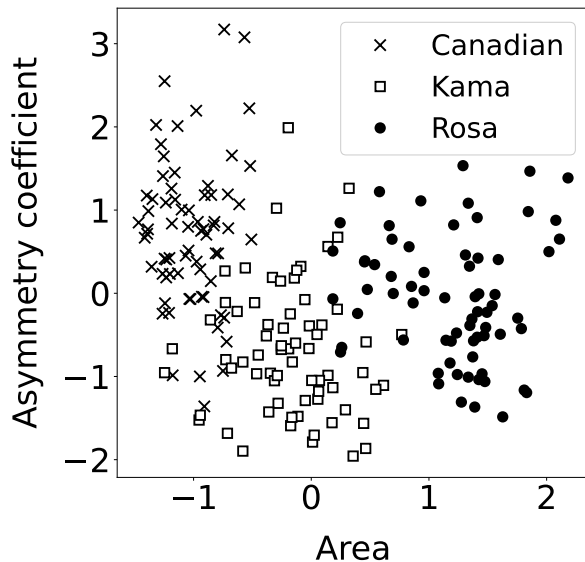
Neural network



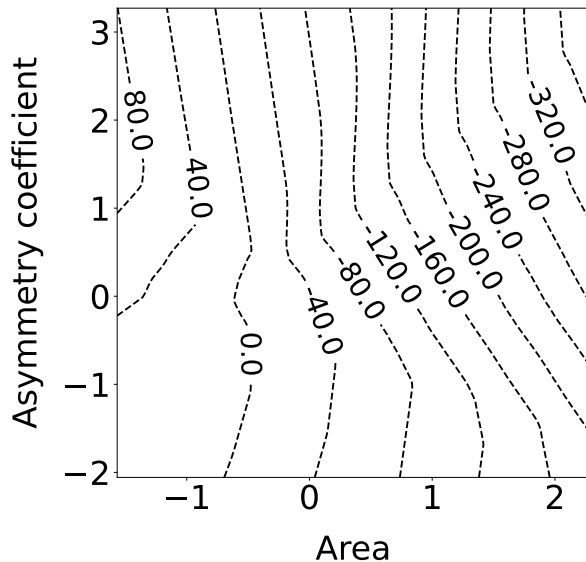
$$P(\tilde{y} = k \mid \tilde{x} = x) \approx p_{\Theta}(x)_k := \frac{\exp(nl_{nnet}(x)[k])}{\sum_{l=1}^c \exp(nl_{nnet}(x)[l])}$$

Non-concave log likelihood maximized via stochastic gradient descent on batches

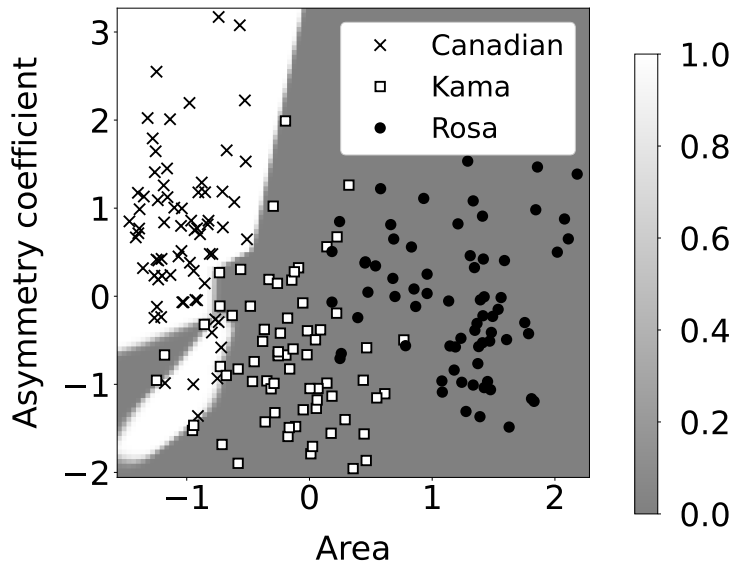
2-layer neural network with 100 hidden variables



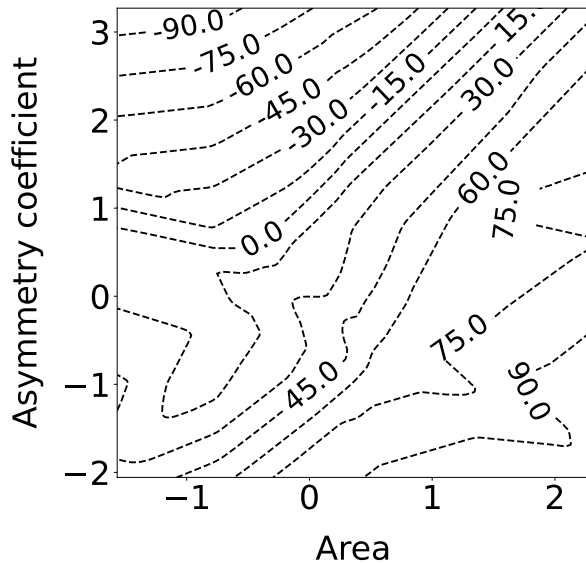
Nonlinear output: Canadian



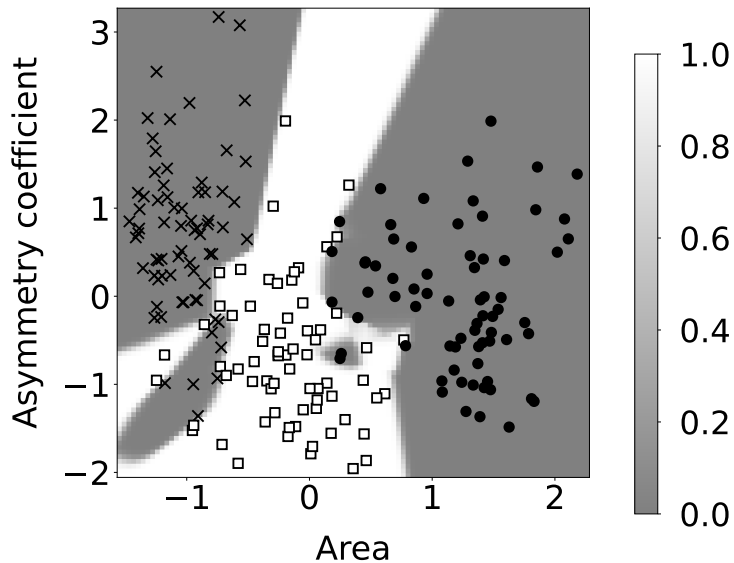
Estimated probability: Canadian



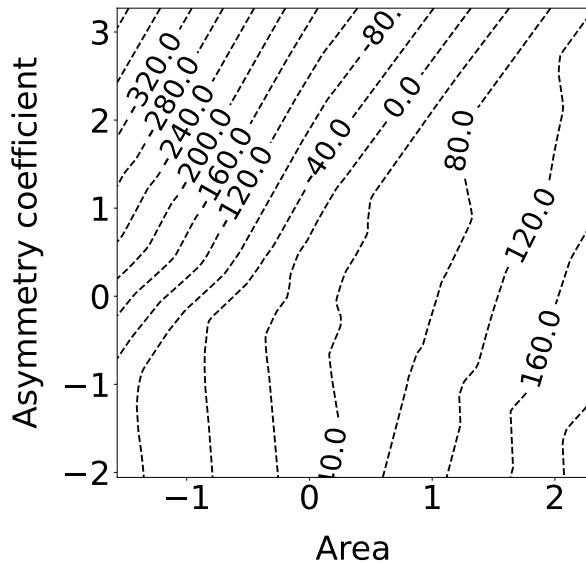
Nonlinear output: Kama



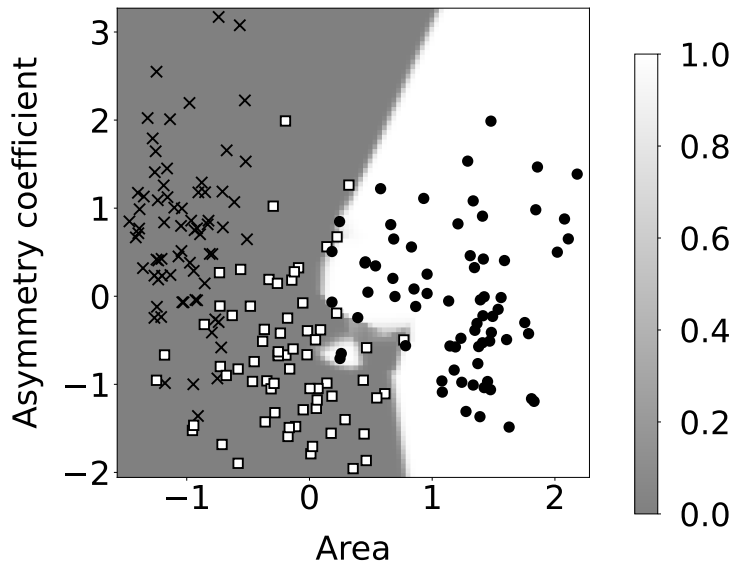
Estimated probability: Kama



Nonlinear output: Rosa



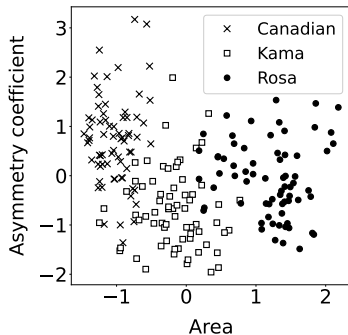
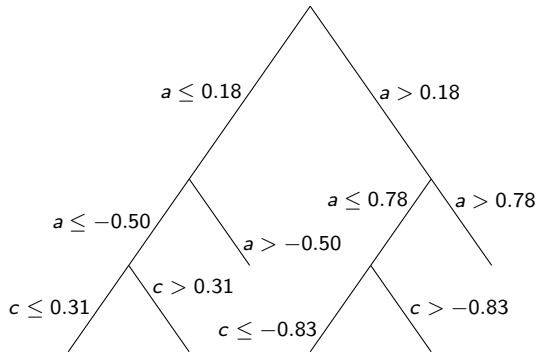
Estimated probability: Rosa



Classification tree

Area (a)

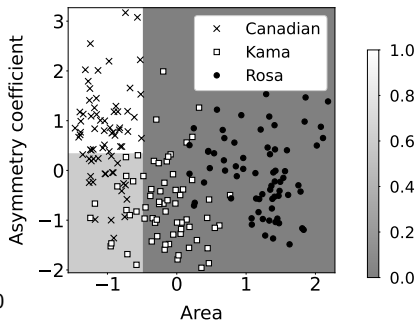
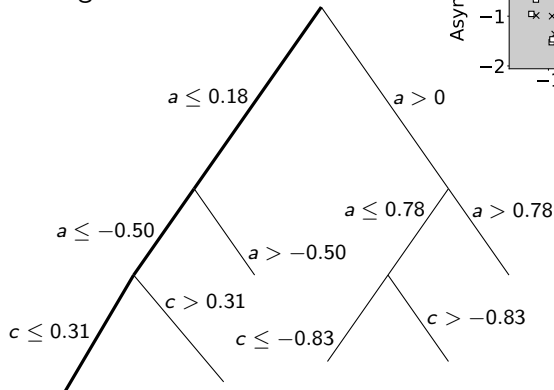
Asymmetric coefficient (c)



$$P(\tilde{y} = \text{Canadian} \mid \tilde{a} = -1, \tilde{c} = -1)$$

Probability estimates assigned to leaves

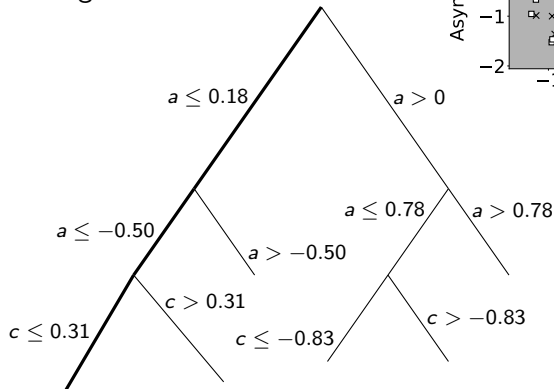
Fraction within region in training set



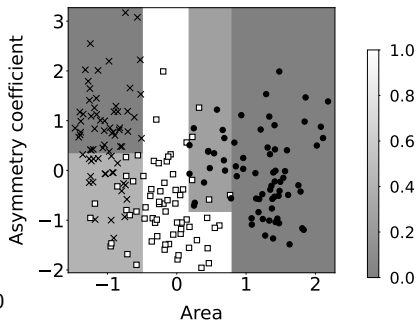
$$P(\tilde{y} = \text{Kama} \mid \tilde{a} = -1, \tilde{c} = -1)$$

Probability estimates assigned to leaves

Fraction within region in training set



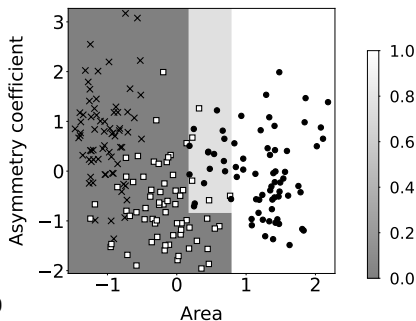
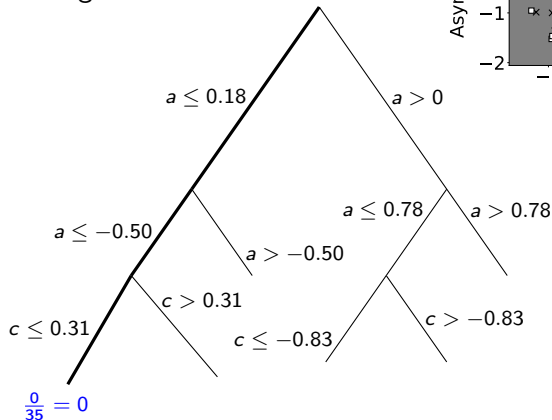
$$\frac{14}{35} = 0.4$$



$$P(\tilde{y} = \text{Rosa} \mid \tilde{a} = -1, \tilde{c} = -1)$$

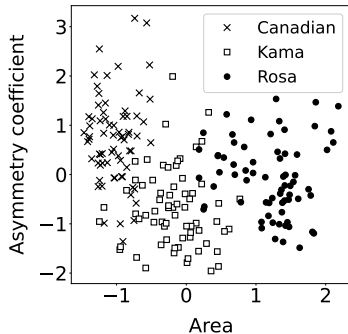
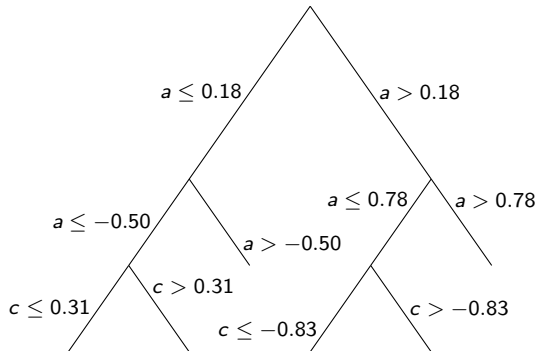
Probability estimates assigned to leaves

Fraction within region in training set



Classification tree

Tree built iteratively, adding bifurcation that most increases the log likelihood



Ensembles

Problem: Simple trees underfit / Complex trees overfit

Solution: Combine multiple simple trees

Three main strategies to obtain individual trees:

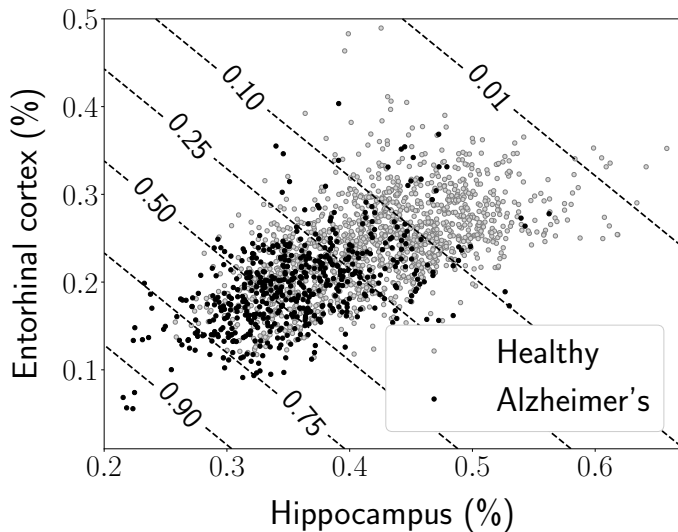
1. **Bagging:** Train on resampled datasets generated via *bootstrapping*
2. **Random forests:** Train *randomized* trees on resampled datasets generated via bootstrapping
3. **Boosting:** Train *complementary* trees that fit residuals of previous trees (scaled down to avoid overfitting)

Generative Models

Discriminative Models

Evaluation

Probability estimates



How should we evaluate them?

Evaluation

We focus on binary classification (2 classes)

- ▶ Threshold probabilities and evaluate binary classification estimates
- ▶ Assess discrimination ability of the probabilities
- ▶ Assess calibration of probabilities

Metrics for binary classification estimates

Accuracy: Fraction of **total** examples that are **correctly classified**

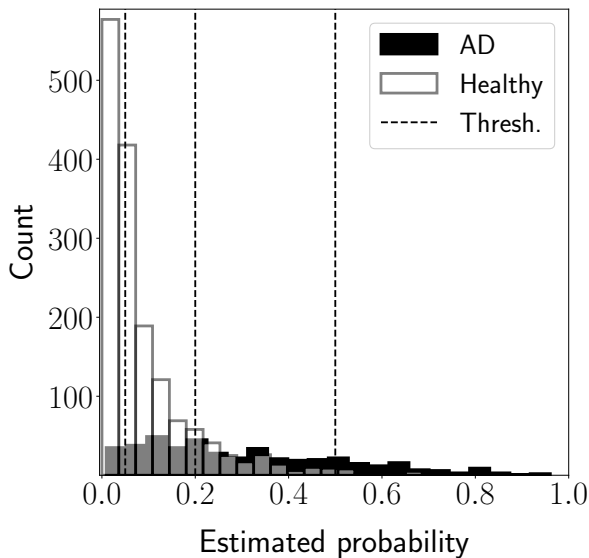
True positive rate (TPR): Fraction of **positive** examples that are correctly classified

False positive rate (FPR): Fraction of **negative** examples that are incorrectly classified

Precision: Fraction of examples **predicted as positive** that are correctly classified

F1 score: Harmonic mean of TPR and precision

Alzheimer's diagnostics



Metrics

Threshold = 0.05

Accuracy = 0.58

TPR = 0.93

FPR = 0.52

Precision = 0.33

F1 score = 0.49

Threshold = 0.2

Accuracy = 0.82

TPR = 0.61

FPR = 0.13

Precision = 0.57

F1 score = 0.59

Threshold = 0.5

Accuracy = 0.81

TPR = 0.19

FPR = 0.01

Precision = 0.78

F1 score = 0.31

TPR - FPR tradeoff

Threshold = 0.05

Accuracy = 0.58

TPR= 0.93

FPR= 0.52

Precision = 0.33

F1 score = 0.49

Threshold = 0.2

Accuracy = 0.82

TPR= 0.61

FPR= 0.13

Precision = 0.57

F1 score = 0.59

Threshold = 0.5

Accuracy = 0.81

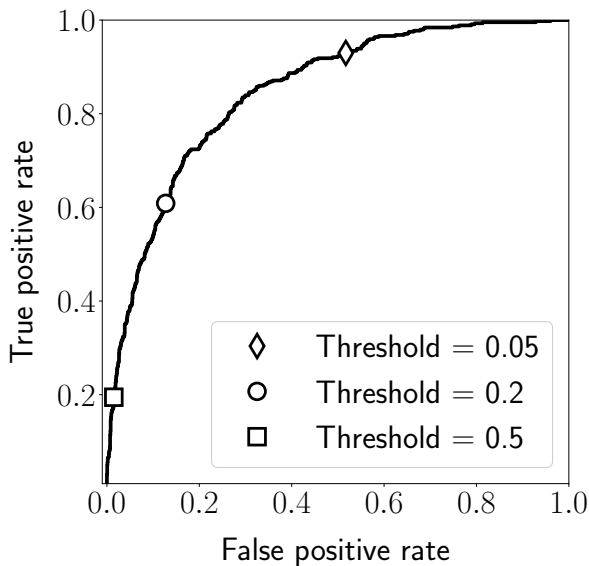
TPR= 0.19

FPR= 0.01

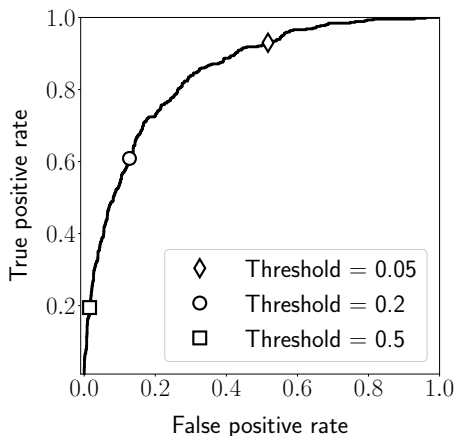
Precision = 0.78

F1 score = 0.31

Receiver operating characteristic (ROC) curve

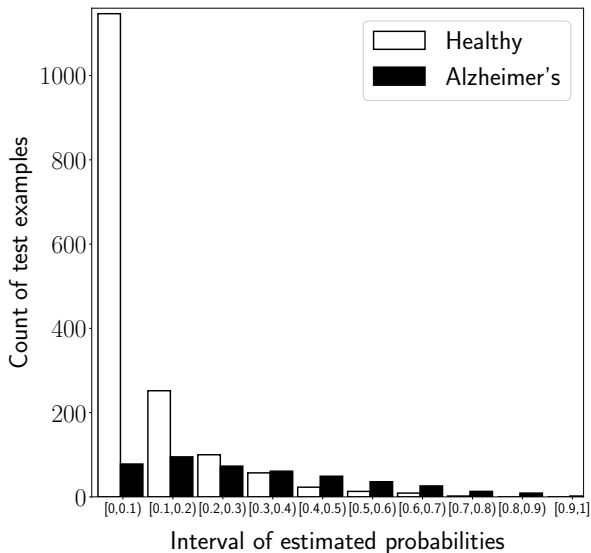


Area under ROC curve (AUROC or AUC) = 0.847

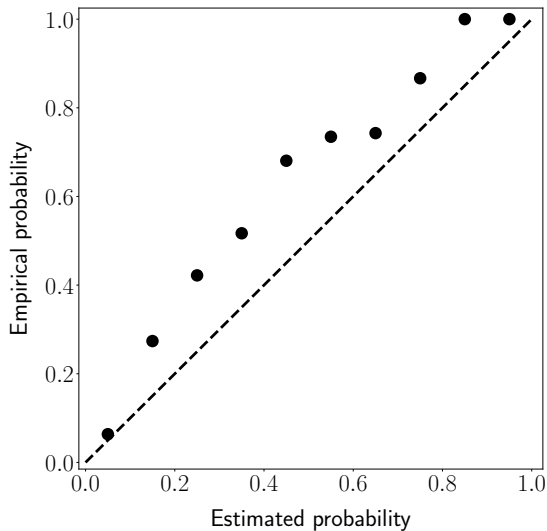


Fraction of negative - positive examples such that estimated probability is higher for positive example

Estimated probabilities vs empirical probabilities



Calibration



Brier score evaluates both discrimination and calibration

What have we learned?

- ▶ Generative models:
 - ▶ Naive Bayes
 - ▶ Gaussian discriminant analysis
- ▶ Discriminative models:
 - ▶ Linear: Logistic regression / Softmax regression
 - ▶ Nonlinear: Neural networks / Classification trees
- ▶ Evaluation