

Logistic Regression

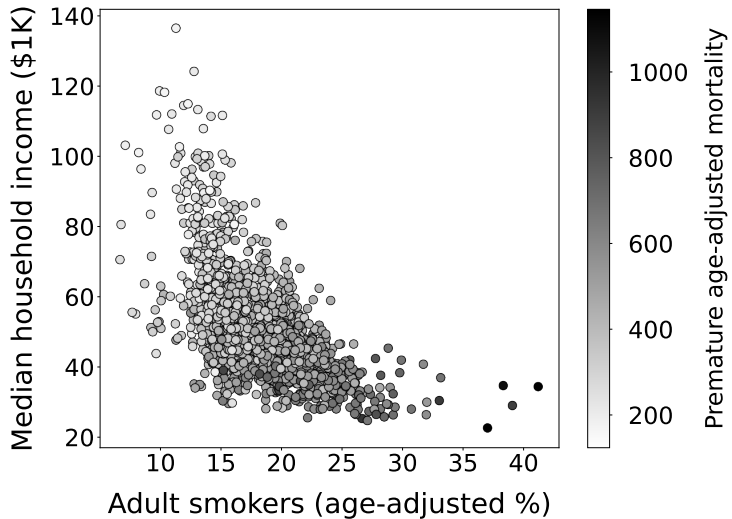
Probability and Statistics for Data Science

Carlos Fernandez-Granda

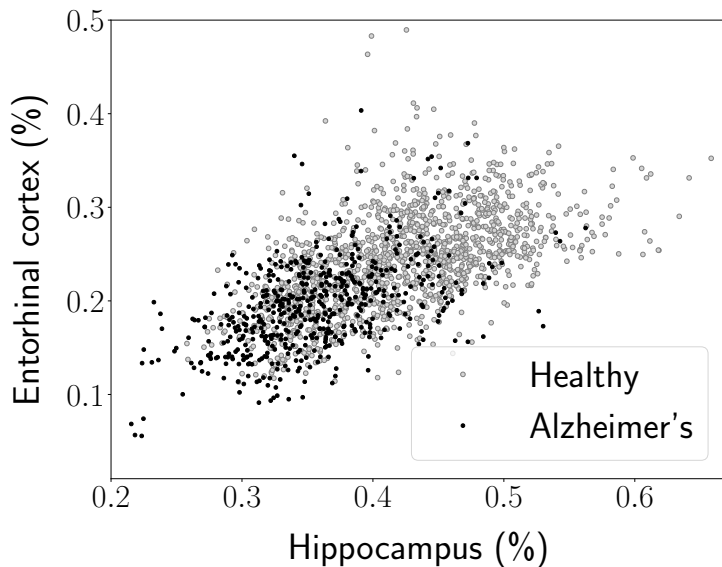


These slides are based on the book [Probability and Statistics for Data Science](#) by Carlos Fernandez-Granda, available for purchase [here](#). A free preprint, videos, code, slides and solutions to exercises are available at <https://www.ps4ds.net>

Regression



Classification



Classification

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

Each **feature** x_i is a d -dimensional vector (e.g. MRI scan)

The label y_i indicates the **class** (e.g. *Alzheimer's* or *healthy*)

Goal: Assign class to new data

Probabilistic modeling

Model features as random vector \tilde{x} and label as random variable \tilde{y}

For new data vector x :

$$\hat{y} := \arg \max_{y \in \{1, 2, \dots, c\}} p_{\tilde{y} | \tilde{x}}(y | x)$$

Is classification easy?

Curse of dimensionality

Unless number of features (entries in \tilde{x}) is very small, it is impossible to estimate $p_{\tilde{y}|\tilde{x}}(y|x)$!

For m binary features we need to estimate 2^m conditional pmfs!

We need assumptions!

Generative vs discriminative approaches

Naive Bayes and Gaussian discriminant analysis are **generative** approaches

First estimate $p_{\tilde{x}|\tilde{y}} / f_{\tilde{x}|\tilde{y}}$ and $p_{\tilde{y}}$, then apply Bayes' rule

Discriminative approaches estimate $p_{\tilde{y}|\tilde{x}}$ directly

Logistic regression

Discriminative *binary* classification method (two classes 0 and 1)

Goal: Use linear model to approximate $p_{\tilde{y}|\tilde{x}}(1|x)$

First idea:

$$p_{\tilde{y}|\tilde{x}}(1|x) = \beta^T x + \alpha$$

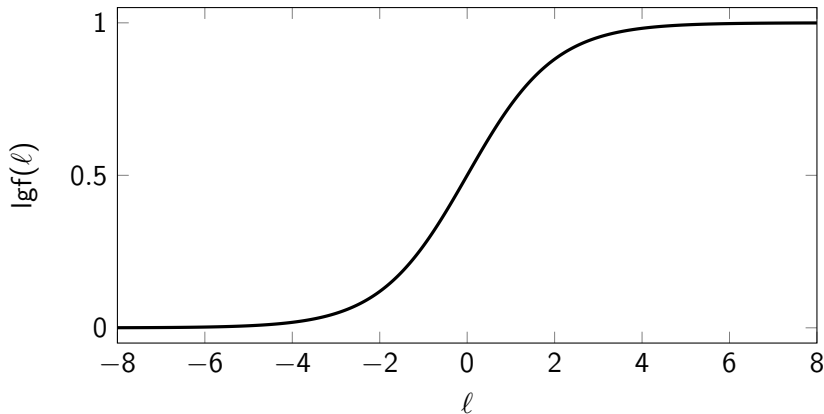
Problem: For most values of x *not a valid probability*

Solution:

Map linear function of features to $[0,1]$ using [link function](#)

Logistic function

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

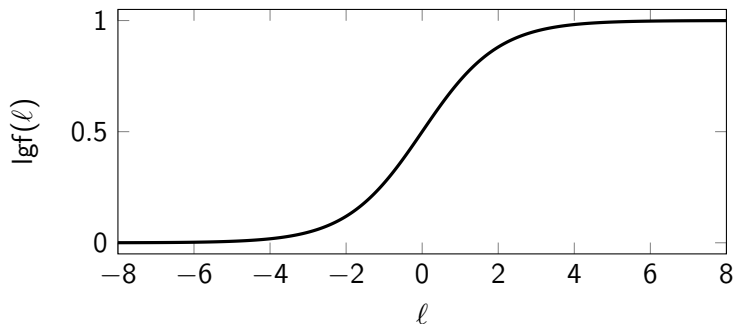


Logistic regression

Generalized linear model

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

Properties of the logistic function



1. Monotone

$$\beta^T \mathbf{x}_1 + \alpha > \beta^T \mathbf{x}_2 + \alpha \implies \text{lgf}(\beta^T \mathbf{x}_1 + \alpha) > \text{lgf}(\beta^T \mathbf{x}_2 + \alpha)$$

2. Continuous

$$\beta^T \mathbf{x}_1 + \alpha \approx \beta^T \mathbf{x}_2 + \alpha \implies \text{lgf}(\beta^T \mathbf{x}_1 + \alpha) \approx \text{lgf}(\beta^T \mathbf{x}_2 + \alpha)$$

3. Saturates: Large negative / positive inputs mapped to 0 / 1

Logits

The logit function is the **inverse** of the logistic function

Inputs to the logistic function are called **logits**

Odds

Ratio between probability of an event and probability of complement

If $P(A) = 0.5$, odds are 1. If $P(A) = 0.75$, odds are 3

For $P(A) := p = \text{lgf}(\ell)$

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

In logistic regression, log odds are affine function of the features

$$\log \left(\frac{P(\tilde{y} = 1 \mid \tilde{x} = x)}{1 - P(\tilde{y} = 1 \mid \tilde{x} = x)} \right) = \beta^T x + \alpha$$

Parameter estimation

How do we estimate the logistic-regression parameters β and α ?

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

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

$$\begin{aligned}\mathcal{L}_{XY}(\alpha, \beta) &:= \mathbb{P}(\tilde{y}_1 = y_1, \dots, \tilde{y}_n = y_n \mid \tilde{x}_1 = x_1, \dots, \tilde{x}_n = x_n) \\ &= \prod_{i=1}^n \mathbb{P}(\tilde{y}_i = y_i \mid \tilde{x}_1 = x_1, \dots, \tilde{x}_n = x_n) \\ &= \prod_{i=1}^n \mathbb{P}(\tilde{y}_i = y_i \mid \tilde{x}_i = x_i)\end{aligned}$$

Likelihood and log-likelihood

$$p_{\alpha,\beta}(x) := \text{lgf}(\beta^T x + \alpha)$$

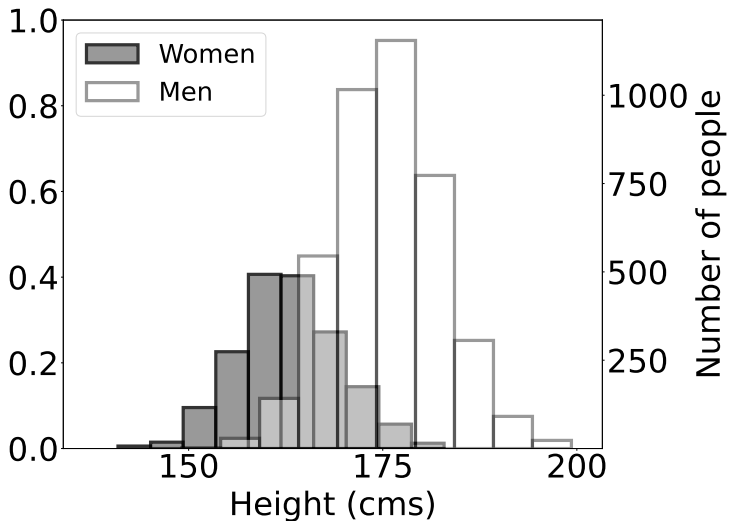
$$\begin{aligned}\mathcal{L}_{XY}(\alpha, \beta) &= \prod_{i=1}^n \text{P}(\tilde{y}_i = y_i \mid \tilde{x}_i = x_i) \\ &= \prod_{\{i: y_i=0\}} (1 - p_{\alpha,\beta}(x_i)) \prod_{\{l: y_l=1\}} p_{\alpha,\beta}(x_l)\end{aligned}$$

$$\log \mathcal{L}_{XY}(\alpha, \beta) = \sum_{\{i: y_i=0\}} \log(1 - p_{\alpha,\beta}(x_i)) + \sum_{\{l: y_l=1\}} \log p_{\alpha,\beta}(x_l)$$

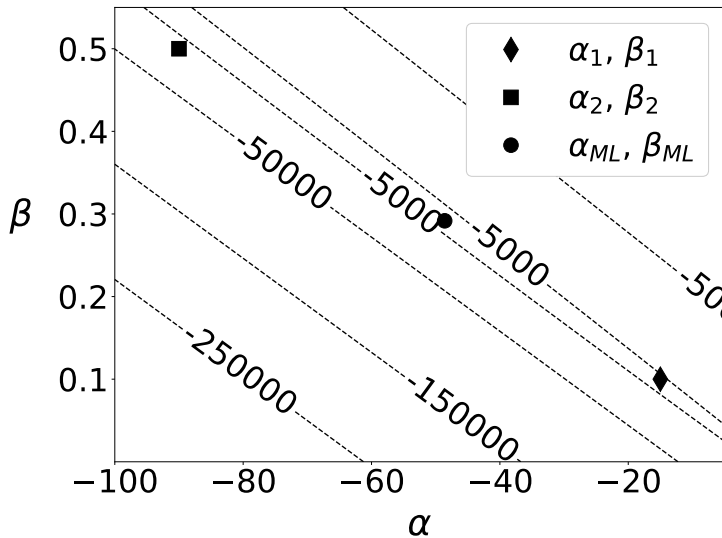
No closed-form solution, but concave

Maximized via iterative methods (gradient ascent, Newton's method, iterative reweighted least squares)

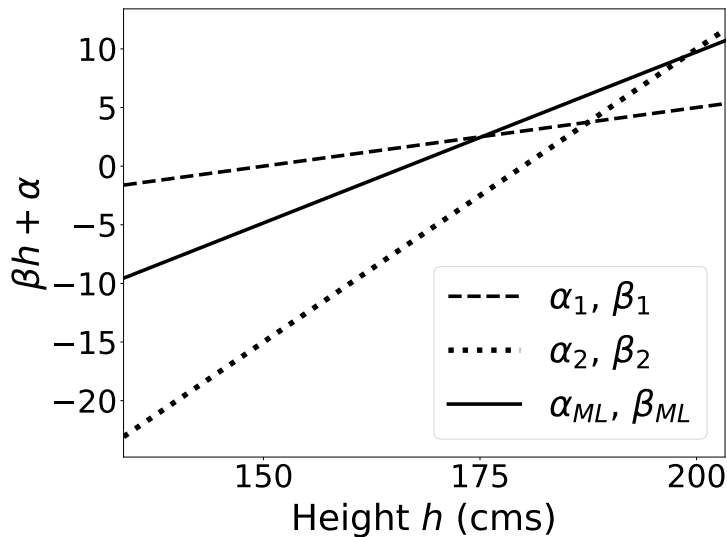
Classification according to height



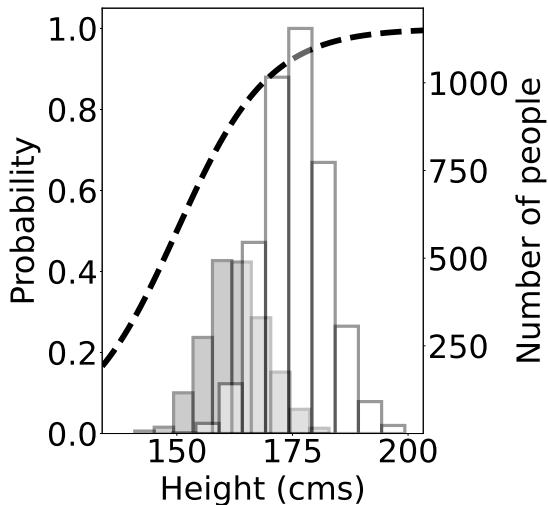
Log-likelihood



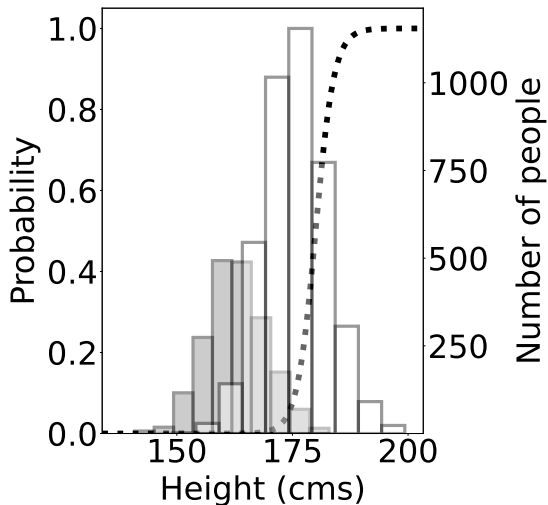
Logits



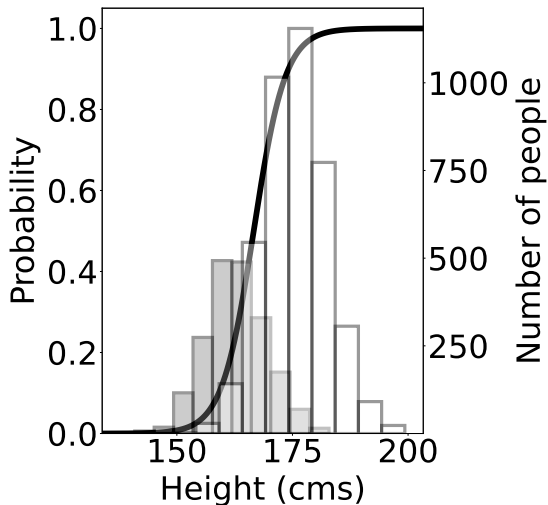
$$\alpha_1 := -15, \beta_1 := 0.1$$



$$\alpha_2 := -90, \beta_2 := 0.5$$



$$\alpha_{\text{ML}} := -48.6, \beta_{\text{ML}} := 0.29$$



Diagnosis of Alzheimer's disease

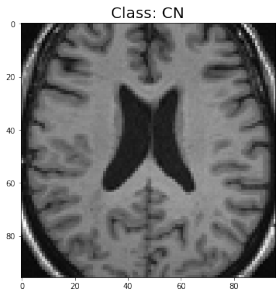
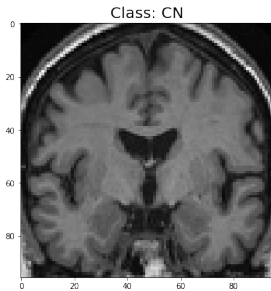
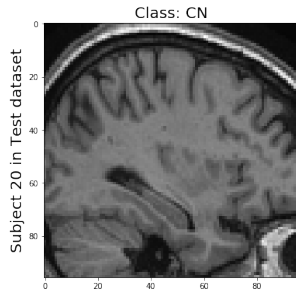
Neurodegenerative disease causing 60 – 70% cases of dementia

Diagnosis via positron-emission tomography is invasive and very costly

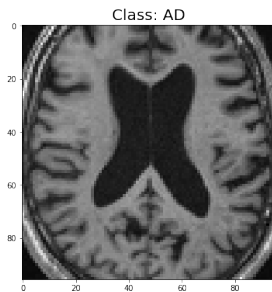
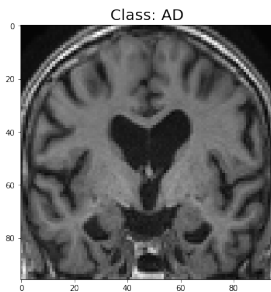
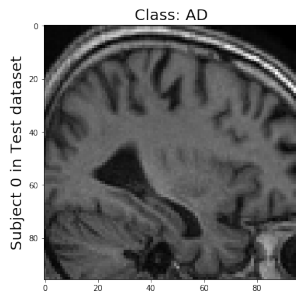
Structural MRI is non-invasive and less costly

Goal: Diagnose Alzheimer's using MRI scans

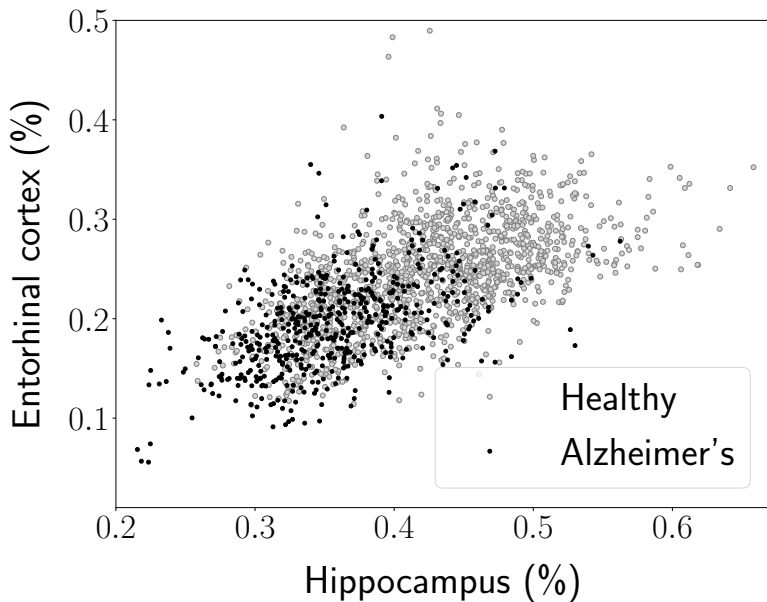
Cognitively-normal patient



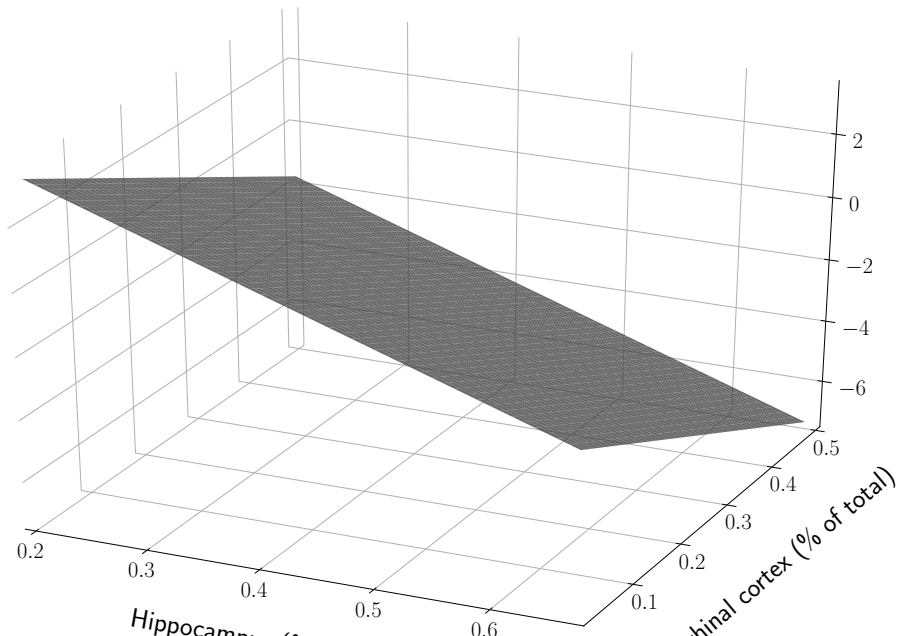
Alzheimer's patient



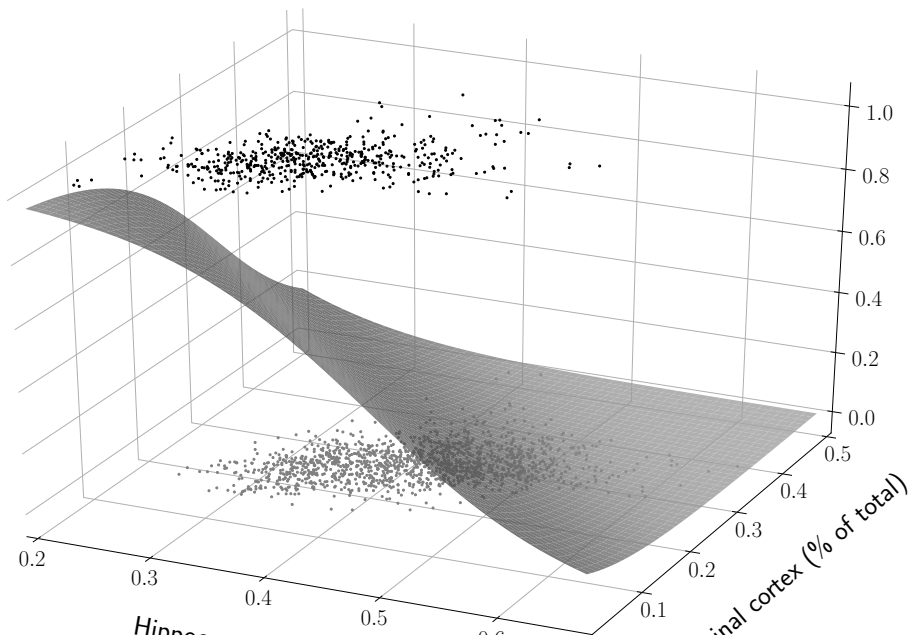
Alzheimer's diagnosis



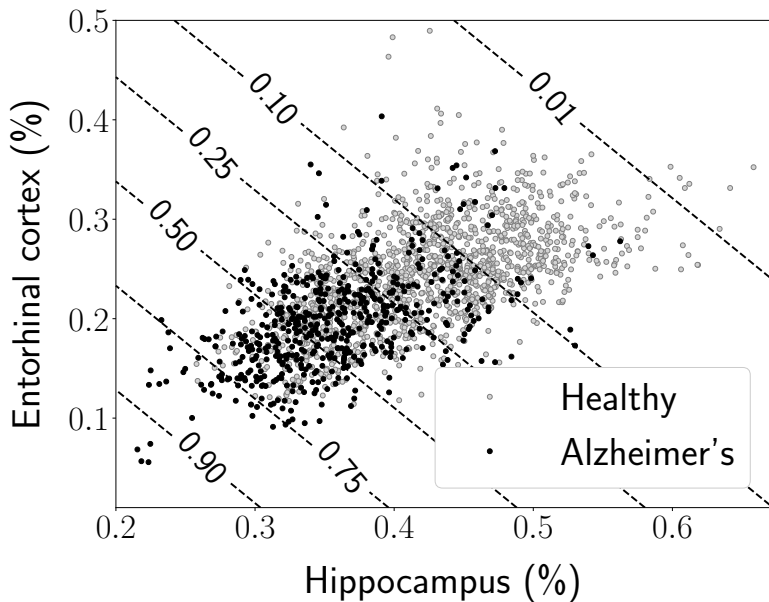
$$-11.9 x_{\text{hippocampus}} - 10.5 x_{\text{entorhinal}} + 5.9$$



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



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



Regularization

Feature collinearity produces noisy coefficients and overfitting (as in linear regression)

Solution: Regularization

Regularized cost function:

$$-\log \mathcal{L}_{XY}(\alpha, \beta) + \lambda \|\beta\|_2^2$$

where λ is a regularization parameter

What have we learned?

How logistic regression works:

- ▶ Link function maps **linear** functions of features (logits) to probability estimates
- ▶ Parameters are obtained by maximizing the **likelihood**