

CS-GY 6923: Lecture 5

Linear Classification, Logistic Regression, Gradient Descent

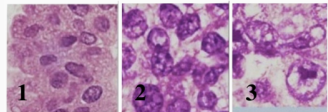
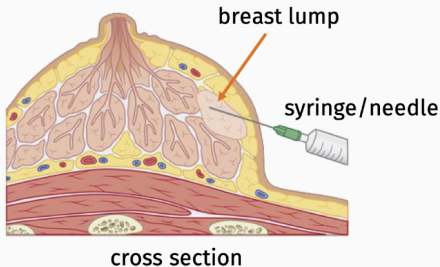
NYU Tandon School of Engineering, Akbar Rafiey

- HW02 is posted. Due on March 05.
- Lab 03 due next week February 26.

Motivating problem

Breast Cancer Biopsy: Determine if a breast lump in a patient is malignant (cancerous) or benign (safe).

- Collect cells from lump using fine needle biopsy.
- Stain and examine cells under microscope.
- Based on certain characteristics (shape, size, cohesion) determine if likely malignant or not).



Motivating problem

Demo: `demo_breast_cancer.ipynb`

Data: UCI machine learning repository

Breast Cancer Wisconsin (Original) Data Set

Download: [Data Folder](#), [Data Set Description](#)

Abstract: Original Wisconsin Breast Cancer Database



Data Set Characteristics:	Multivariate	Number of Instances:	699	Area:	Life
Attribute Characteristics:	Integer	Number of Attributes:	10	Date Donated	1992-07-15
Associated Tasks:	Classification	Missing Values?	Yes	Number of Web Hits:	564320

[https://archive.ics.uci.edu/ml/datasets/breast+cancer+wisconsin+\(original\)](https://archive.ics.uci.edu/ml/datasets/breast+cancer+wisconsin+(original))

Features: 9 numerical scores about cell characteristics (Clump Thickness, Uniformity, Marginal Adhesion, etc.)

Motivating problem

Data: $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$.

$\mathbf{x}_i = [1, 5, 4, \dots, 2]$ contains score values.

Label $y_i \in \{0, 1\}$ is 0 if benign cells, 1 if malignant cells.

Goal: Based on scores predict if a sample of cells is malignant or benign.

Approach:

- Naive Bayes Classifier can be extended to \mathbf{x} with numerical values (instead of binary values as seen before).

What are other classification algorithms people have heard of?

k-nearest neighbor method

***k*-NN algorithm:** a simple but powerful baseline for classification.

Training data: $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$ where $y_1, \dots, y_n \in \{1, \dots, q\}$.

Classification algorithm:

Given new input \mathbf{x}_{new} ,

- Compute $sim(\mathbf{x}_{new}, \mathbf{x}_1), \dots, sim(\mathbf{x}_{new}, \mathbf{x}_n)$.¹
- Let $\mathbf{x}_{j_1}, \dots, \mathbf{x}_{j_k}$ be the training data vectors with highest similarity to \mathbf{x}_{new} .
- Predict y_{new} as *majority*(y_{j_1}, \dots, y_{j_k}).

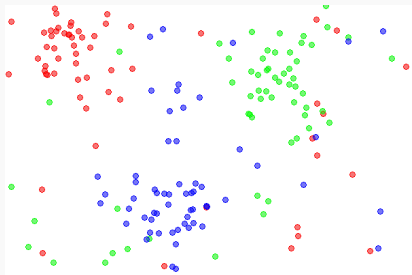
k=5

2 → 1

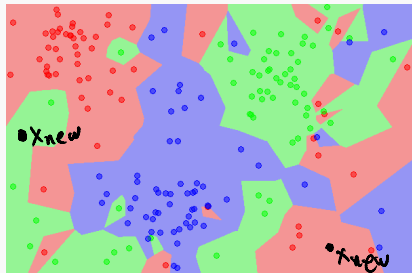
3 → 0

¹ $sim(\mathbf{x}_{new}, \mathbf{x}_i)$ is any chosen similarity function, like $1 - \|\mathbf{x}_{new} - \mathbf{x}_i\|_2$.

k -nearest neighbor method



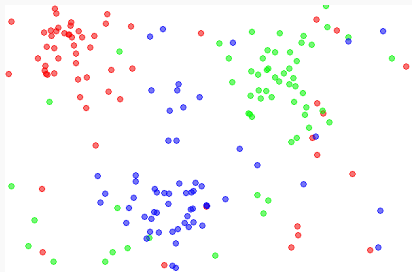
Data



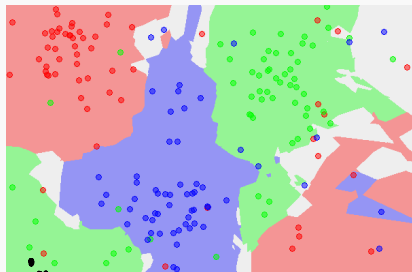
1-NN classifier

- Smaller k , more complex classification function.
- Larger k , more robust to noisy labels.

k -nearest neighbor method



Data



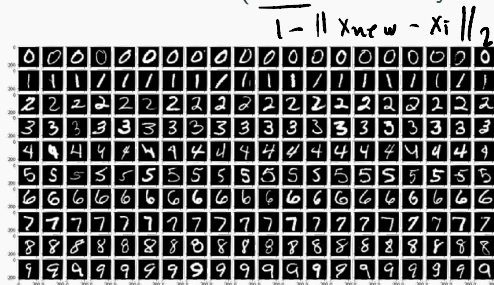
x_{new}
↳ green 5-NN classifier

- Smaller k , more complex classification function.
- Larger k , more robust to noisy labels.
- Works remarkably well for many datasets.

MNIST image data

Especially good for large datasets with lots of repetition.

Example: works well on MNIST (≈ 95% accuracy out-of-the-box.):



$k=3$

Can be improved to 99.5% with a fancy similarity function!²

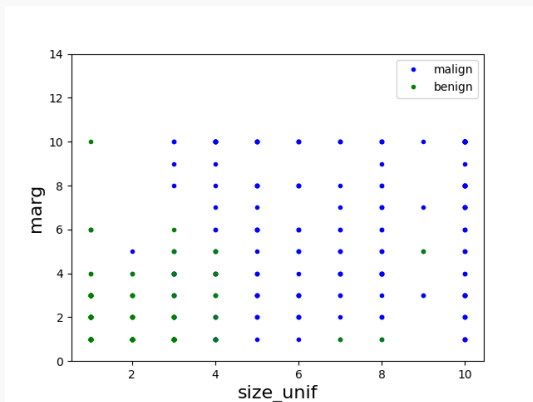
One issue is that prediction can be computationally intensive...

²We will revisit this when we talk about kernel methods.

Linear classification

Begin by plotting data

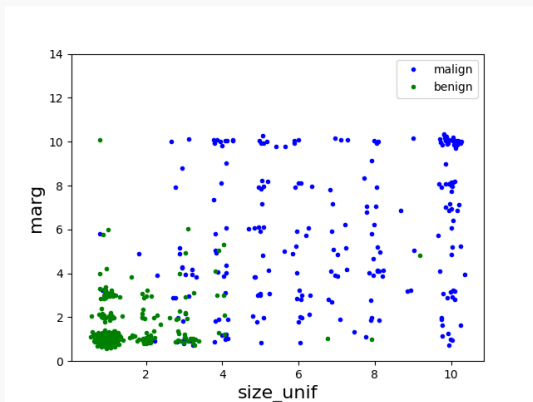
We pick two variables, Margin Adhesion and Size Uniformity and plot a scatter plot. Points with label 1 (malignant) are plotted in blue, those with label 2 (benign) are plotted in green.



Lots of overlapping points! Hard to get a sense of the data.

Plotting with jitter

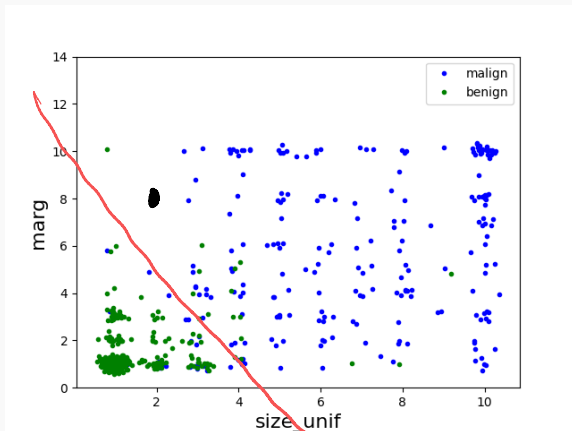
Simple + Useful Trick: data jittering. Add tiny random noise (using e.g. `np.random.randn`) to data to prevent overlap.



Noise is only for plotting. It is not added to the data for training, testing, etc.

Brainstorming

Any ideas for possible classification rules for this data?



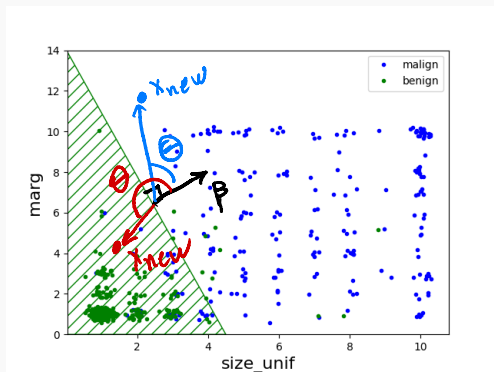
Linear classifier

Given vector of predictors $\mathbf{x}_i \in \mathbb{R}^d$ (here $d = 2$) find a parameter vector $\beta \in \mathbb{R}^d$ and threshold λ .

- Predict $y_i = 0$ if $\langle \mathbf{x}_i, \beta \rangle \leq \lambda$.
- Predict $y_i = 1$ if $\langle \mathbf{x}_i, \beta \rangle > \lambda$



$$\langle a, b \rangle = \|a\| \cdot \|b\| \cdot \cos(\theta)$$



Line has equation $\langle \mathbf{x}, \beta \rangle = \lambda$.

Linear classifier

As long as we append a 1 onto each data vector \mathbf{x}_i (i.e. a column of ones onto the data matrix \mathbf{X}) like we did for linear regression, an equivalent function is:

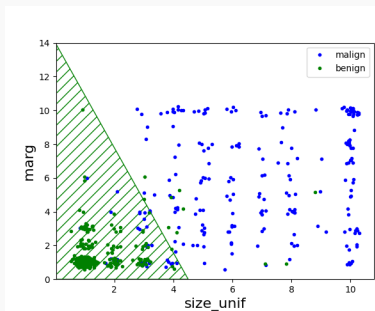
- Predict $y_i = 0$ if $\langle \mathbf{x}_i, \boldsymbol{\beta} \rangle \leq 0$.
- Predict $y_i = 1$ if $\langle \mathbf{x}_i, \boldsymbol{\beta} \rangle > 0$

$$[1 \ x_1 \ x_2]$$

$$[\beta_0 \ \beta_1 \ \beta_2]$$

$$\beta_0 + \beta_1 x_1 + \beta_2 x_2 \leq 0$$

$$\mathbf{X} = \begin{bmatrix} x_1 & x_{11} & x_{12} \\ x_2 & x_{12} & x_{22} \\ \vdots & \vdots & \vdots \end{bmatrix}$$



Line has equation $\langle \mathbf{x}, \boldsymbol{\beta} \rangle = 0$.

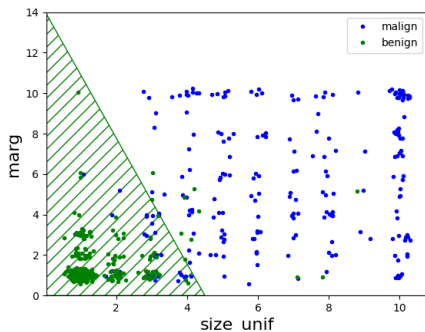
$$\mathbf{x}_{new}$$

$$\mathbf{x}_{new} = [1 \ x_1 \ x_2]$$

Linear classification

Standard approach for binary classification of real-valued data:

- Find parameter vector β .
- For input data vector \mathbf{x} , predict 0 if $\beta^T \mathbf{x} \leq \lambda$ and 1 if $\beta^T \mathbf{x} > \lambda$ for some threshold λ .³



³Can always assume $\lambda = 0$ if \mathbf{x} has an intercept term.

Question: How do we find a good linear classifier automatically?

Loss minimization approach (first attempt):

- **Model⁴:**

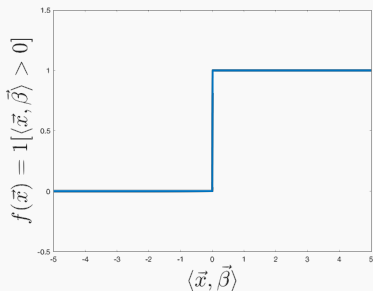
$$f_{\beta}(\mathbf{x}) = \mathbb{1} [\langle \mathbf{x}, \beta \rangle > 0]$$

- **Loss function:** “0 - 1 Loss”

$$L(\beta) = \sum_{i=1}^n |f_{\beta}(\mathbf{x}_i) - y_i|$$

⁴ $\mathbb{1}[\text{event}]$ is the indicator function: it evaluates to 1 if the argument inside is true, 0 if false.

Problem with 0 – 1 loss:



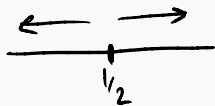
- The loss function $L(\beta)$ is not differentiable because $f_{\beta}(\mathbf{x})$ is discontinuous.
- Impossible to take the gradient, very hard to minimize loss to find optimal β .
- Non-convex function (will make more sense next lecture).

Linear classifier via square loss

Loss minimization approach (second attempt):

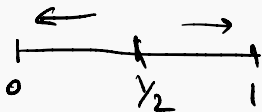
- Model:

$$f_{\beta}(\mathbf{x}) = \mathbb{1} [\langle \mathbf{x}, \beta \rangle > 1/2]$$



- Loss function: "Square Loss"

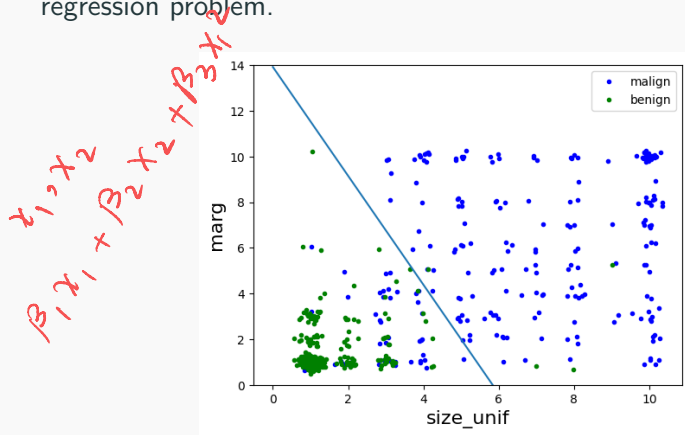
$$L(\beta) = \sum_{i=1}^n (\langle \mathbf{x}_i, \beta \rangle - y_i)^2$$



Intuitively tries to make $\langle \mathbf{x}, \beta \rangle$ close to 0 for examples in class 0, close to 1 for examples in class 1.

Linear classifier via square loss

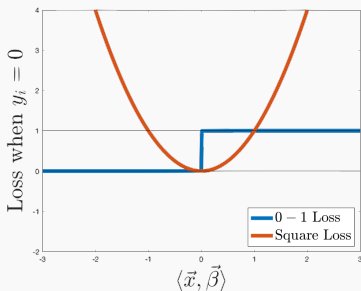
We can solve for β by just solving a least squares multiple linear regression problem.



Do you see any issues here?

Linear classifier via square loss

Problem with square loss:



- Loss increases if $\langle \mathbf{x}, \boldsymbol{\beta} \rangle > 1$ even if correct label is 1. Or if $\langle \mathbf{x}, \boldsymbol{\beta} \rangle < 0$ even if correct label is 0.
- Intuitively we don't want to "punish" these cases.

Logistic regression

Let $h_{\beta}(x)$ be the **logistic function**:

$$h_{\beta}(x) = \frac{1}{1 + e^{-\langle \beta, x \rangle}}$$

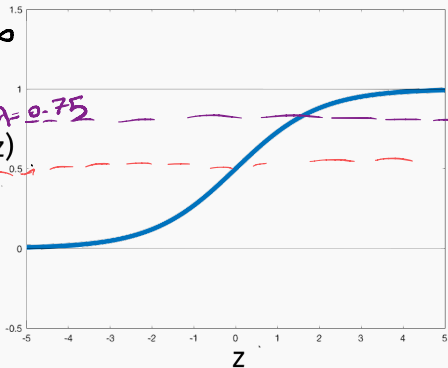
$$\text{if } \langle \beta, x \rangle = 0 \\ \Rightarrow h_{\beta}(x) = \frac{1}{1 + e^0}$$

$$= \frac{1}{1+1} = \frac{1}{2}$$

if $\langle x, \beta \rangle \rightarrow -\infty$

$$\Rightarrow h_{\beta}(x) = \frac{1}{1 + e^{\infty}}$$

$= 0$



$$\text{if } \langle \beta, x \rangle \rightarrow +\infty \\ \Rightarrow h_{\beta}(x) = \frac{1}{1 + e^{-\infty}}$$

$$= \frac{1}{1+0} = 1$$

Logistic regression

Loss minimization approach (this works!):

- **Model:** Let $h_{\beta}(\mathbf{x}) = \frac{1}{1+e^{-\langle \beta, \mathbf{x} \rangle}}$

$$\begin{aligned}f_{\beta}(\mathbf{x}) &= \mathbb{1} [h_{\beta}(\mathbf{x}) > 1/2] \\ &= \mathbb{1} [\langle \mathbf{x}, \beta \rangle > 0]\end{aligned}$$

- **Loss function:** “Logistic loss” aka “binary cross-entropy loss”

$$L(\beta) = - \sum_{i=1}^n y_i \log(h_{\beta}(\mathbf{x}_i)) + (1 - y_i) \log(1 - h_{\beta}(\mathbf{x}_i))$$

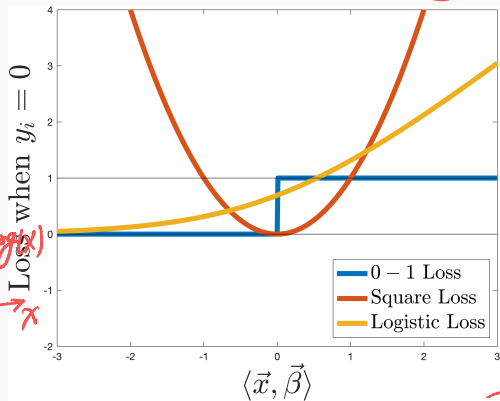
if $y_i = 0 \Rightarrow$

if $y_i = 1$

Logistic loss

Logistic Loss:

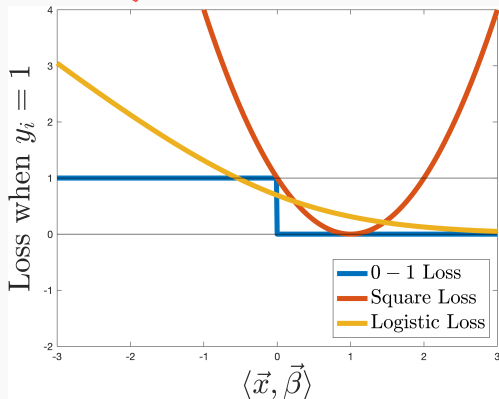
$$L(\beta) = - \sum_{i=1}^n y_i \log(h_{\beta}(\mathbf{x})) + (1 - y_i) \log(1 - h_{\beta}(\mathbf{x}))$$



Logistic loss

Logistic Loss:

$$L(\beta) = - \sum_{i=1}^n y_i \log(h_{\beta}(\mathbf{x})) + (1 - y_i) \log(1 - h_{\beta}(\mathbf{x}))$$



Loss minimization approach:

- Given training data $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$ where $\mathbf{x}_i \in \mathbb{R}^d$ and $y_i \in \{0, 1\}$.
- Minimize “Logistic loss” aka “binary cross-entropy loss”

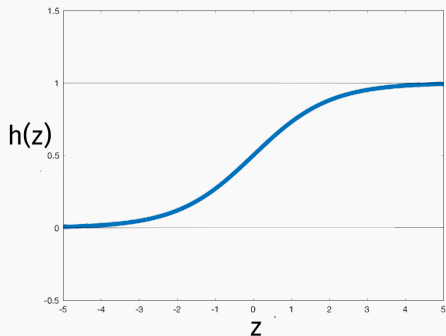
$$L(\beta) = - \sum_{i=1}^n y_i \log(h(\beta^T \mathbf{x}_i)) + (1 - y_i) \log(1 - h(\beta^T \mathbf{x}_i))$$

- Above $h(z)$ is the logistic/sigmoid function: $h(z) = \frac{1}{1+e^{-z}}$

Prediction: Predict $y_i = 1$ if $\beta^T \mathbf{x}_i \geq 0$, predict 0 otherwise.

Logistic regression

Let $h(z)$ be the logistic/sigmoid function: $h(z) = \frac{1}{1+e^{-z}}$



Can think of this function as mapping $\mathbf{x}^T \boldsymbol{\beta}$ to a probability that the true label is 1. If $\mathbf{x}^T \boldsymbol{\beta} \gg 0$ then the probability is close to 1, if $\mathbf{x}^T \boldsymbol{\beta} \ll 0$ then the probability is close to 0.

Exercise

Why not minimize:

$$L(\beta) = \sum_{i=1}^n \left(y_i - h(\mathbf{x}^T \beta) \right)^2 ?$$

Exercise

Why not minimize:

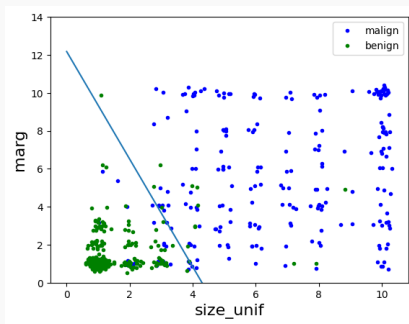
$$L(\beta) = \sum_{i=1}^n \left(y_i - h(\mathbf{x}^T \beta) \right)^2 ?$$

Answer: This is actually a pretty reasonable thing to do. An important issue however is that the loss here is not convex, which makes it hard to find the β that minimizes the loss.

Log-loss on the other hand is convex. More on this later.

Logistic loss

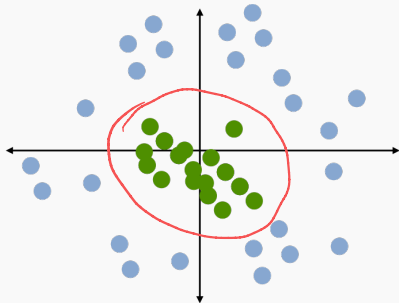
- Convex function in β , can be minimized using gradient descent.
- Works well in practice.
- Good Bayesian motivation.
- Easily combined with non-linear data transformations.



Fit using logistic regression/log loss.

Non-linear transformations

How would we learn a classifier for this data using logistic regression?

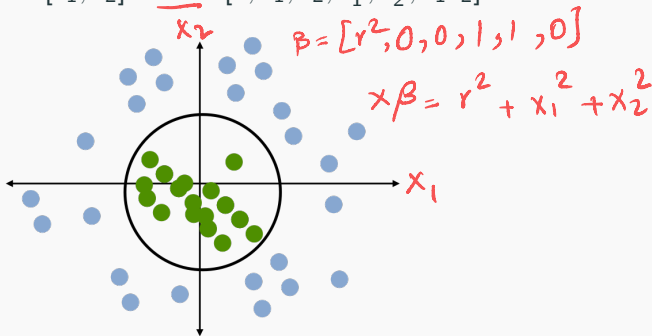


This data is not linearly separable or even approximately linearly separable.

Non-linear transformations

Transform each $\mathbf{x} = [x_1, x_2]$ to $\mathbf{x} = [1, x_1, x_2, x_1^2, x_2^2, x_1x_2]$

$$x_1^2 + x_2^2 = r^2$$

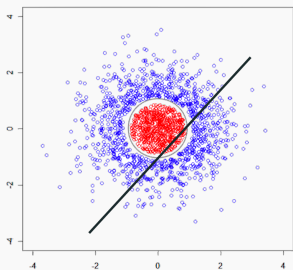


- Predict class 1 if $x_1^2 + x_2^2 < \lambda$.
- Predict class 0 if $x_1^2 + x_2^2 \geq \lambda$.

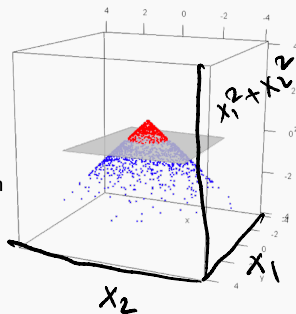
This is a linear classifier on our transformed data set. Logistic regression might learn $\beta = [r^2, 0, 0, 1, 1, 0]$.

Non-linear transformations

View as mapping data to a higher dimensional space, where it is linearly separable.



feature
transformation



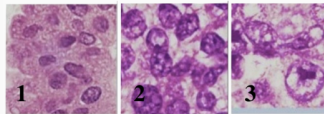
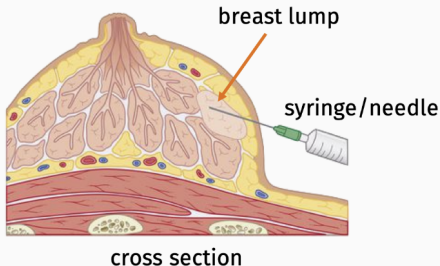
Lots more on this in future lecture!

Error in classification

Once we have a classification algorithm, how do we judge its performance?

- **Simplest answer:** Error rate = fraction of data examples misclassified in test set.
- What are some issues with this approach?

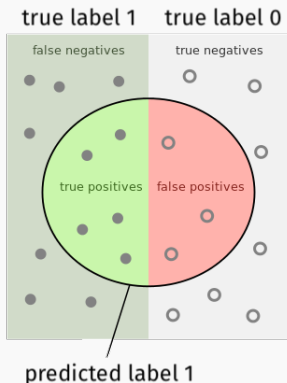
Think back to motivating problem of breast cancer detection.



Error in classification

- **Precision:** Fraction of positively labeled examples (label 1) which are correct.
- **Recall:** Fraction of true positives that we labeled correctly with label 1.

Question: Which should we optimize for medical diagnosis?

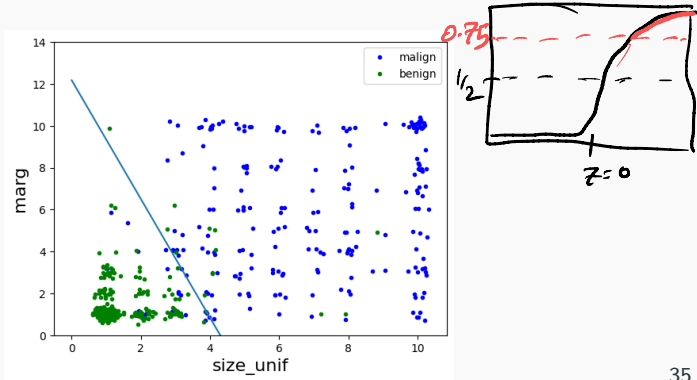


$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$
$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

Error in classification

Possible logistic regression workflow:

- Learn $\vec{\beta}$ and compute $h_{\vec{\beta}}(\vec{x}_i) = \frac{1}{1+e^{-\langle \vec{x}_i, \vec{\beta} \rangle}}$ for all \vec{x}_i .
- Predict $y_i = 0$ if $h_{\vec{\beta}}(\vec{x}_i) \leq \lambda$, $y_i = 1$ if $h_{\vec{\beta}}(\vec{x}_i) > \lambda$.
- Default value of λ is $1/2$. How does changing λ affect precision and recall ?



Possible logistic regression workflow:

- Learn $\vec{\beta}$ and compute $h_{\vec{\beta}}(\vec{x}_i) = \frac{1}{1+e^{-\langle \vec{x}_i, \vec{\beta} \rangle}}$ for all \vec{x}_i .
- Predict $y_i = 0$ if $h_{\vec{\beta}}(\vec{x}_i) \leq \lambda$, $y_i = 1$ if $h_{\vec{\beta}}(\vec{x}_i) > \lambda$.
- Default value of λ is $1/2$. How does changing λ affect precision and recall ?

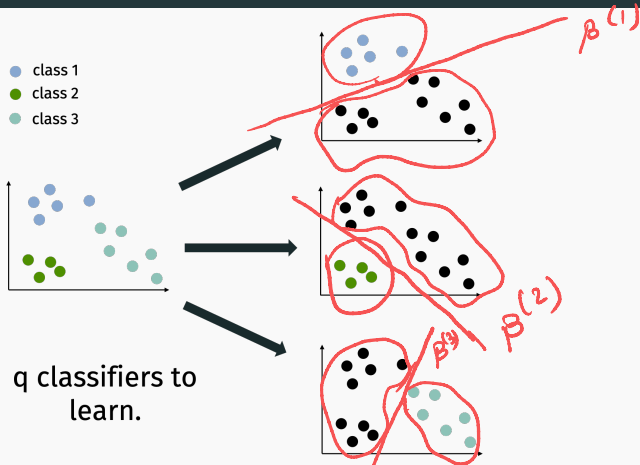
This is very heuristic. There are other methods for handling “class imbalance” which can often lead to good overall error, but poor precision or recall. Techniques include weighting the loss function to care more about false negatives, or subsampling the larger class.

What about when $y \in \{1, \dots, q\}$ instead of $y \in \{0, 1\}$?

Two common options for reducing multi-class problems to binary problems:

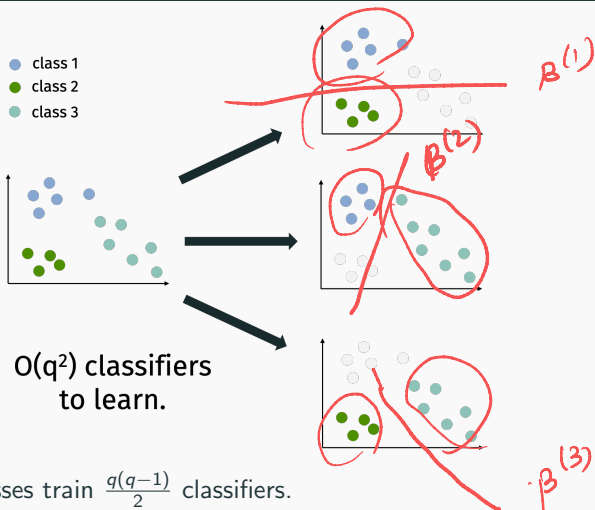
- One-vs.-all (most common, also called one-vs.-rest)
- One-vs.-one (slower, but can be more effective)

One vs. rest



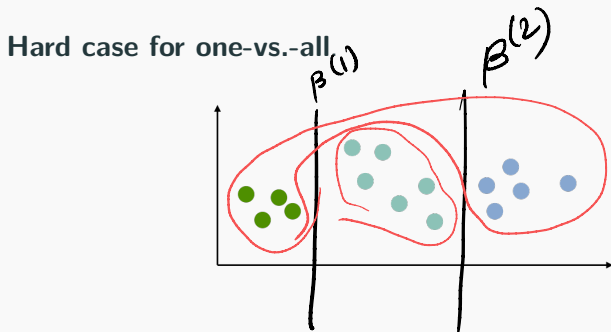
- For q classes train q classifiers. Obtain parameters $\beta^{(1)}, \dots, \beta^{(q)}$.
- Assign y to class i if $\langle \beta^{(i)}, \mathbf{x}_{new} \rangle \geq 0$. Could be ambiguous!
- **Better:** Assign y to class i with maximum value of $h(\langle \beta^{(i)}, \mathbf{x}_{new} \rangle)$.

One vs. one



- For q classes train $\frac{q(q-1)}{2}$ classifiers.
- Assign y to class i which wins in the most number of head-to-head comparisons.

One vs. one



- One-vs.-one would be a better choice here.
- Also tends to work better when there is class imbalance.
- But one-vs.-one can be super expensive! E.g when $q = 100$ or $q = 1000$.

Multiclass logistic regression

More common modern alternative: If we have q classes, train a single model with q parameter vectors $\beta^{(1)}, \dots, \beta^{(q)}$, and predict class $i = \arg \max_j \langle \beta^{(j)}, \mathbf{x} \rangle$.

Same idea as one-vs.-rest, but we treat $[\beta^{(1)}, \dots, \beta^{(q)}]$ as a single length qd parameter vector which we optimize to minimize a single joint loss function. We do not train the parameter vectors separately.

What's a good loss function?

Multiclass logistic regression

Softmax function:

$$\begin{bmatrix} \langle \beta^{(1)}, \mathbf{x} \rangle \\ \vdots \\ \langle \beta^{(q)}, \mathbf{x} \rangle \end{bmatrix} \xrightarrow{\text{softmax}} \begin{bmatrix} e^{\langle \beta^{(1)}, \mathbf{x} \rangle} / \sum_{i=1}^q e^{\langle \beta^{(i)}, \mathbf{x} \rangle} \geq 0 \\ \vdots \\ e^{\langle \beta^{(q)}, \mathbf{x} \rangle} / \sum_{i=1}^q e^{\langle \beta^{(i)}, \mathbf{x} \rangle} \geq 0 \end{bmatrix}$$

= |

Softmax takes in a vector of numbers and converts it to a vector of probabilities:

$$\begin{bmatrix} -10 & 4 & 1 & 0 & -5 \end{bmatrix} \rightarrow \begin{bmatrix} .00 & \underline{.94} & .04 & .02 & .00 \end{bmatrix}$$

$$\frac{e^4}{(e^{-10} + e^4 + e^1 + e^0 + e^{-5})} = 0.94$$

Multiclass logistic regression

Multi-class cross-entropy:

$$\begin{aligned}L(\beta^{(1)}, \dots, \beta^{(q)}) &= - \sum_{i: y_i=1} \log \frac{e^{\langle \beta^{(1)}, \mathbf{x}_i \rangle}}{\sum_{j=1}^q e^{\langle \beta^{(j)}, \mathbf{x}_i \rangle}} - \dots - \sum_{i: y_i=q} \log \frac{e^{\langle \beta^{(q)}, \mathbf{x}_i \rangle}}{\sum_{j=1}^q e^{\langle \beta^{(j)}, \mathbf{x}_i \rangle}} \\ &= - \sum_{i=1}^n \sum_{\ell=1}^q \underbrace{\mathbb{1}[y_i = \ell]}_{0 \text{ or } 1} \cdot \log \frac{e^{\langle \beta^{(\ell)}, \mathbf{x}_i \rangle}}{\sum_{j=1}^q e^{\langle \beta^{(j)}, \mathbf{x}_i \rangle}}\end{aligned}$$

Binary cross-entropy:

$$\begin{aligned}L(\beta) &= - \sum_{i=1}^n y_i \log(h(\beta^T \mathbf{x}_i)) + (1 - y_i) \log(1 - h(\beta^T \mathbf{x}_i)) \\ &= - \sum_{i: y_i=1} \log(h(\beta^T \mathbf{x}_i)) - \sum_{i: y_i=0} \log(1 - h(\beta^T \mathbf{x}_i))\end{aligned}$$

Not exactly the same... but can show equivalent if you set $\beta^{(0)} = \beta$ and $\beta^{(1)} = -\beta$.

Error in (multiclass) classification

Confusion matrix for k classes:

Pred->	1	2	...	K
Real↓				
1	1	0.5	0	
2	<u>0.5</u>	1		0.99
...			1	
K				1

- Entry i, j is the fraction of class i items classified as class j .
- Useful to see whole matrix to visualize where errors occur.

Optimization

Goal: Minimize the logistic loss:

$$L(\beta) = - \sum_{i=1}^n y_i \log(h(\beta^T \mathbf{x}_i)) + (1 - y_i) \log(1 - h(\beta^T \mathbf{x}_i))$$

I.e. find $\beta^* = \arg \min L(\beta)$. How should we do this?

Logistic regression gradient

$$L(\beta) = - \sum_{i=1}^n y_i \log(h(\beta^T \mathbf{x}_i)) + (1 - y_i) \log(1 - h(\beta^T \mathbf{x}_i))$$

Let $\mathbf{X} \in \mathbb{R}^{d \times n}$ be our data matrix with $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ as rows.

Let $\mathbf{y} = [y_1, \dots, y_n]$. A calculation gives (verify!):

$$\left[\begin{array}{c} \frac{\partial L(\beta)}{\partial \beta_1} \\ \vdots \end{array} \right] \leftarrow \nabla L(\beta) = \underbrace{\mathbf{X}^T}_{d \times n} \underbrace{(h(\mathbf{X}\beta) - \mathbf{y})}_{n \times 1}$$

$\hookrightarrow n \times 1$

where $h(\mathbf{X}\beta) = \frac{1}{1+e^{-\mathbf{X}\beta}}$. Here all operations are entrywise. I.e in Python you would compute:

1	<code>h = 1 / (1 + np.exp(-X@beta))</code>
2	<code>grad = np.transpose(X)@(h - y)</code>

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

$$\begin{array}{ccc} h(\mathbf{X}\beta) & & \\ \swarrow \downarrow & & \searrow \\ n \times d & \xrightarrow{d \times 1} & n \times 1 \end{array}$$

Logistic regression gradient

To find β minimizing $L(\beta)$ we typically start by finding a β where:

$$\nabla L(\beta) = \mathbf{X}^T (h(\mathbf{X}\beta) - \mathbf{y}) = \mathbf{0}$$
$$\mathbf{X}^T (\mathbf{X}\beta - \mathbf{y}) = \mathbf{0} \Rightarrow \mathbf{X}^T \mathbf{X} \beta = \mathbf{X}^T \mathbf{y} \rightarrow \beta = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- In contrast to what we saw when minimizing the squared loss for linear regression, there's no simple closed form expression for such a β !
- This is the typical situation when minimizing loss in machine learning: linear regression was a lucky exception.
- **Main question:** How do we minimize a loss function $L(\beta)$ when we can't explicitly compute where its gradient is $\mathbf{0}$?

Minimizing loss functions

Always an option: Brute-force search. Test as many possible values for β and just see which gives the smallest value of $L(\beta)$.

- As we saw on Lab 1, this actually works okay for low-dimensional problems (e.g. when β has 1 or 2 entries).
- **Problem:** Super computationally expensive in high-dimension. For $\beta \in \mathbb{R}^d$, run time grows as:

$$\begin{array}{ccccccc} [\beta_1 & \beta_2 & \dots & \beta_d] & & & \\ \downarrow & \downarrow & & \downarrow & & & \\ 100 & 100 & & 100 & & & 100^d \end{array}$$

Minimizing loss functions

Much Better idea. Some sort of guided search for a good of β .

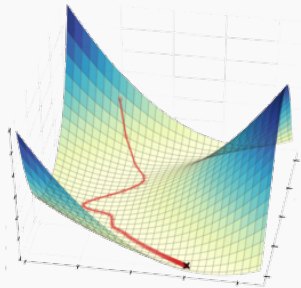
- Start with some $\beta^{(0)}$, and at each step try to change β slightly to reduce $L(\beta)$.
- Hopefully find an approximate minimizer for $L(\beta)$ much more quickly than brute-force search.
- **Concrete goal:** Find β with

$$L(\beta) < \min_{\beta} L(\beta) + \epsilon$$

for some small error term ϵ .

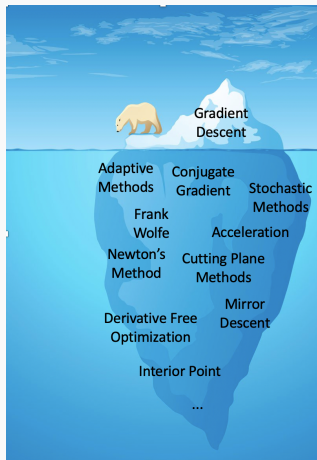
Gradient descent

Gradient descent: A greedy search algorithm for minimizing functions of multiple variables (including loss functions) that often works amazingly well.



The single most important computational tool in machine learning.
And it's remarkable simple + easy to implement.

Optimization algorithms



Just one method in a huge class of algorithms for numerical optimization. All of these methods are important in ML.