

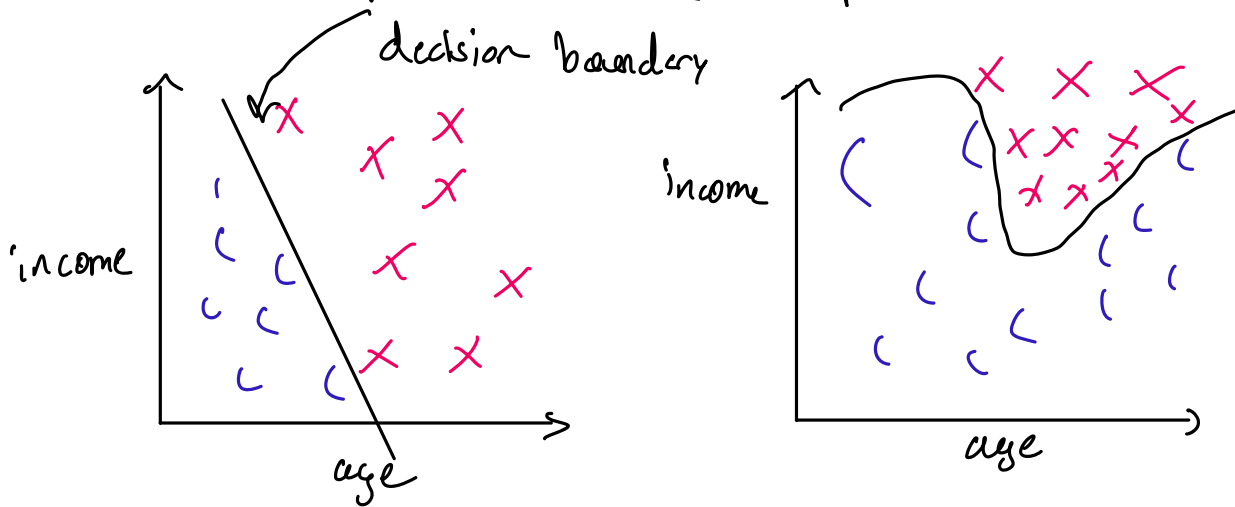
Classifiers

- Given a sample of n observations w/ d features/predictors each.
- Some observations belong to class C ; others do not

example: bank loans = observations
features = income & age
classes = [defaulted, not defaulted]

goal: predict whether someone will fault or not based on income & age

- represent each object in d dimensional space
↳ sample point / feature vector / independent variables



decision boundary: separates classes

Overfitting: when decision boundaries become "snake-like"
↳ real test: will boundary classify test data well?

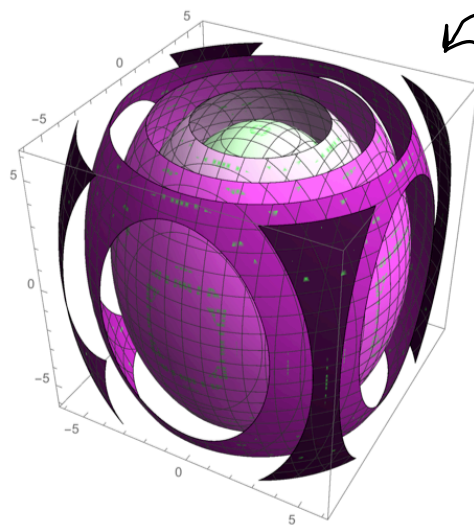
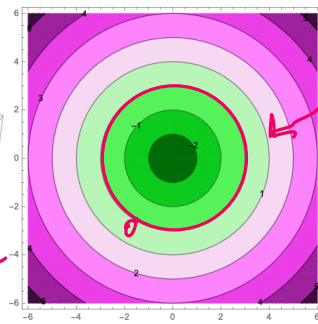
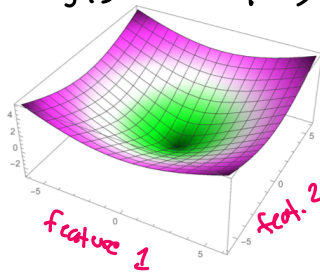
aka "predictor functions" or "discriminants"

Some classifiers compute a **decision function**: function $f(x)$ that maps sample point \vec{x} to scalar s.t. $f(x)$ is positive if sample is in class C and vice versa

For such classifiers, decision boundary is defined by $\{x \in \mathbb{R}^d, f(x)=0\}$ infinite set of points

$\{x: f(x)=0\}$ is also called an **isosurface** of $f(\cdot)$ for **isovalue** 0
* f has other isosurfaces & isovalues e.g. $\{x: f(x)=1\}$

$$f(x,y) = \sqrt{x^2 + y^2} - 3$$



2D isoc contours in
3D space

linear classifier: decision boundary is linear (line/plane/hyperplane...)
 ↳ most use linear decision function
 SVM's don't have linear

Math Review

Vectors: $x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{bmatrix} = [x_1 \ x_2 \ x_3 \ x_4 \ x_5]^T$
 ↳ point in SD space

Conventions: uppercase roman: Matrix, Random Variable, Set
 lowercase roman: Vector
 greek: Scalar

Other scalars: $n = \# \text{ sample points}$, $d = \# \text{ features/dimension of sample points}$
 $i, j, k = \text{integer indices}$

Inner product: $X \cdot Y = x_1 y_1 + \dots + x_d y_d$
 also written in matrix notation: $x^T y$

clearly, $f(x) = w \cdot x + d$ is a linear function (linear in x)

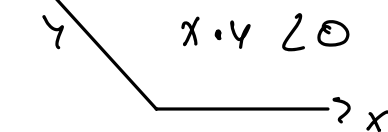
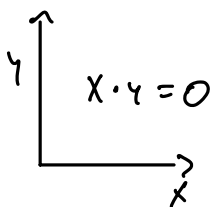
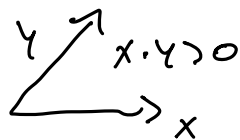
Euclidean Norm: $\|x\| = \sqrt{x^T x} = \sqrt{x_1^2 + x_2^2 + \dots + x_d^2}$

$\|x\|$ is length of \vec{x} (euclidean length)

Given \vec{x} , $\frac{\vec{x}}{\|\vec{x}\|}$ is a **unit vector** (length = 1)

normalization →

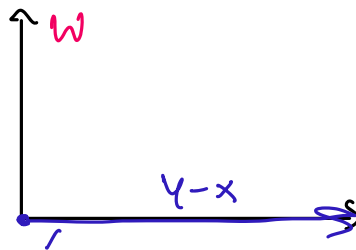
Use dot products to compute angles: $\cos \theta = \frac{x \cdot y}{\|x\| \cdot \|y\|} = \frac{x}{\|x\|} \cdot \frac{y}{\|y\|}$



given $f(x) = w \cdot x + d$, decision boundary is

Set $H = \{x : w \cdot x = -d\}$
Set H is a **hyperplane** (line in 2D, plane in 3D)

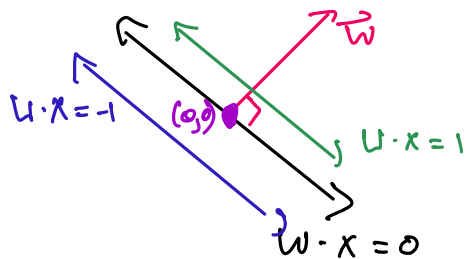
- dimension $d-1$
- flat
- infinite
- cuts d -space into 2



Theorem: Let x, y be 2 points on H . Then $w \cdot (y-x) = 0$

Proof: $w(y-x) = -d + d = 0 \Rightarrow w \perp x \forall x \in H$

→ w is called **normal vector** of H (perpendicular)



If \vec{w} is a **unit vector**, then $w \cdot x + d$ is the **Signed distance** from point x to hyperplane H i.e. positive on w 's side of H
negative on opposite side of w

Moreover, distance from H to origin is d

⇒ $d = 0$ iff H passes through origin

The coefficients in w , plus d , are the **weights** (or parameters) of the classifier.
or regression coefficients

Input data is **linearly separable** if there exists a hyperplane that separates all samples by their class
→ If not separable, accuracy won't be 100%.

e.g. Simple Classifier: **Centroid Method**

- 1) - compute mean μ_c of all points in class C
- compute μ_x of all points not in class C

2) decision function $f(x) = (\underbrace{\mu_c - \mu_x}_{\text{normal vector}}) \cdot \vec{x} - (\mu_c - \mu_x) \cdot \frac{\mu_c + \mu_x}{2}$

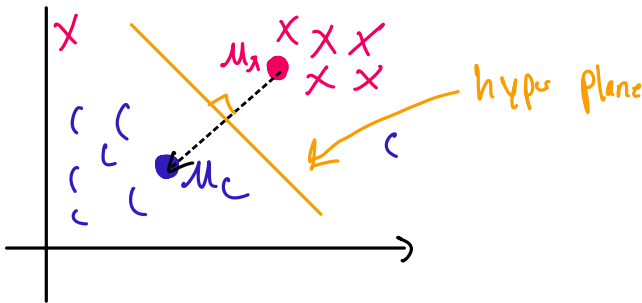
midpoint of mean vectors

→ decision boundary passes through midpoint between μ_c & μ_x

decision boundary is H-plane that bisects line segment w/ endpoints μ_c, μ_x

- when $\vec{x} = \frac{\mu_c + \mu_x}{2}$ $f(x) = 0$

→ decision boundary passes through midpoint



Perceptron algorithm (Frank Rosenblatt, 1957)

- slow
- correct for linearly separable points
- uses a **numerical optimization algorithm**, namely **GRADIENT DESCENT**

Consider n training points $x_1, x_2 \dots x_n$

each training point x_i has a **label** y_i
each row is a datapoint

$$\text{label} = \begin{cases} 1 & x_i \text{ in class } C \\ -1 & \text{e/w} \end{cases}$$

For simplicity, consider only boundaries that pass through origin

goal: find weights \vec{w} s.t. $\vec{x}_i \cdot \vec{w} \geq 0$ if $y_i = 1$

$$\vec{x}_i \cdot \vec{w} < 0 \quad \text{if } y_i = -1$$

equivalently: $y_i x_i \cdot \vec{w} \geq 0$ ← constraint

Idea: define risk function R that is positive if some constraints are violated, 0 otherwise

Then, use optimization algo that minimizes R

Loss Function:

classifiers prediction

$$L(z, y_i) = \begin{cases} 0 & \text{if } y_i z \geq 0 \\ -y_i z & \text{otherwise} \end{cases}$$

← scalar, same sign

Risk Function:
(cost function)
(obj. function)

$$R(w) = \frac{1}{n} \sum_{i=1}^n L(x_i \cdot w, y_i)$$

$$= \frac{1}{n} \sum_{i \in V} -y_i (x_i \cdot w), \quad V \text{ is set of misclassified training points}$$

$$R(x) = \begin{cases} 0 & \text{all training points classified correctly} \\ \text{positive} & \text{O/W} \end{cases}$$

goal: Solve this optimization problem

Find w that minimizes $R(w)$

$R(w)$

