CS 1675: Intro to Machine Learning Intro to Classification (Nearest Neighbors, Logistic Regression, Perceptron)

> Prof. Adriana Kovashka University of Pittsburgh September 27, 2018

Classification

- Given features **x**, predict categorical output y
- For example:
 - Given attributes of a house (e.g. square footage and age built), predict whether it will be bought for the asking price or for less
 - Given temperature, predict whether it will rain, snow, or be sunny
- The rest of the course will cover different supervised approaches to classification

Plan for this lecture

- The simplest classifier: K-Nearest Neighbors
 - Algorithm and example use
 - Generalizing: Distance metrics, weighing neighbors
 - Problems: curse of dimensionality, picking K
- Logistic regression
 - Probability: review
 - Linear regression for classification?
 - Maximum likelihood solution for *logistic* regression
 - Related algorithm: perceptron

Nearest Neighbors: Key Idea

- A type of supervised learning: We want to learn to predict, for a new data point x, its label y (e.g. spam / not spam)
- Don't learn an explicit function F: X \rightarrow Y
- Keep all training data {X, Y}
- For a test example x, find the training example x_i closest to it (e.g. using Euclidean distance)
- Then copy the target label y_i as the label for x

Related Methods / Synonyms

- Instance-based methods
- Exemplar methods
- Memory-based methods
- Non-parametric methods

Instance/Memory-based Learning

Four things make a memory based learner:

• A distance metric

• How many nearby neighbors to look at?

• A weighting function (optional)

• *How to fit with the local points?*

1-Nearest Neighbor Classifier

Four things make a memory based learner:

- A distance metric
 - Euclidean (and others)
- How many nearby neighbors to look at?
 1
- A weighting function (optional)
 Not used
- How to fit with the local points?
 Predict the same output as the nearest neighbor

1-Nearest Neighbor Classifier



f(x) = label of the training example nearest to x

Adapted from Lana Lazebnik

K-Nearest Neighbor Classifier

Four things make a memory based learner:

- A distance metric
 - Euclidean (and others)
- How many nearby neighbors to look at?
 K
- A weighting function (optional)
 Not used
- How to fit with the local points?
 Predict the average output among the nearest neighbors

K-Nearest Neighbor Classifier

- For a new point, find the k closest points from training data (e.g. k=5)
- Labels of the *k* points "vote" to classify



1-nearest neighbor



3-nearest neighbor



5-nearest neighbor



What are the tradeoffs of having a too large k? Too small k?

Formal Definition

- Let x be our test data point, and N_κ(x) be the indices of the k nearest neighbors of x
- Classification:

$$y = \operatorname{argmax}_{c} \#(y_{i} = c)$$
$$y = \operatorname{argmax}_{c} \sum_{i \in N_{K}(x)} I(y_{i} = c)$$

• Regression:

$$y = \frac{1}{K} \sum_{i \in N_K(x)} y_i$$

Example: Predict where this picture was taken



Hays and Efros, IM2GPS: Estimating Geographic Information from a Single Image, CVPR 2008

Example: Predict where this picture was taken



Example: Predict where this picture was taken



6+ million geotagged photos by 109,788 photographers



Scene Matches













Hays and Efros, IM2GPS: Estimating Geographic Information from a Single Image, CVPR 2008

Scene Matches





Croatia



england









France







Barcelona



Paris

Malta

Italy









Austria

Hays and Efros, IM2GPS: Estimating Geographic Information from a Single Image, CVPR 2008





Latvia

Italy



europe







Scene Matches







The Importance of Data



k-Nearest Neighbor

Four things make a memory based learner:

- A distance metric
 - Euclidean (and others)
- How many nearby neighbors to look at?
 k
- A weighting function (optional)
 Not used
- *How to fit with the local points?*
 - Just predict the average output among the nearest neighbors

Distances

- Suppose I want to charge my overall distance more for differences in x₂ direction as opposed to x₁ direction
- Setup A: equal weighing on all directions
- Setup B: more weight on x₂ direction
- Will my neighborhoods be longer in the x₁ or x₂ direction?

Voronoi partitioning

- Nearest neighbor regions
- All points in a region are closer to the seed in that region than to any other seed (black dots = seeds)



Multivariate distance metrics

Suppose the input vectors $\mathbf{x}_1, \mathbf{x}_2, ... \mathbf{x}_N$ are two dimensional:

 $\mathbf{x}_{1} = (x_{1}^{1}, x_{1}^{2}), \mathbf{x}_{2} = (x_{2}^{1}, x_{2}^{2}), \dots, \mathbf{x}_{N} = (x_{N}^{1}, x_{N}^{2}).$



The relative scalings in the distance metric affect region shapes

Adapted from Carlos Guestrin

Distance metrics

- Euclidean:
- Minkowski:

• Mahalanobis:

$$d(\mathbf{x}, \mathbf{z}) = \left[\sum_{\substack{i=1\\D}}^{D} (x_i - z_i)^2\right]^{\frac{1}{2}}$$
$$d(\mathbf{x}, \mathbf{z}) = \left[\sum_{i=1}^{D} |x_i - z_i|^P\right]^{\frac{1}{P}}$$
$$d(\mathbf{x}, \mathbf{z}) = \sqrt{\sum_{i=1}^{D} \frac{(x_i - z_i)^2}{\sigma_i^2}}$$
$$d(\mathbf{x}, \mathbf{z}) = (\mathbf{x} - \mathbf{z})^{\mathrm{T}} A(\mathbf{x} - \mathbf{z})^{\mathrm{T}}$$

(where A is a **positive semidefinite** matrix, i.e. symmetric matrix with all non-negative eigenvalues)

Distance metrics Euclidean



Another generalization: Weighted K-NNs

- Neighbors weighted differently:
 - Use all samples, i.e. K = N
 - Weight on i-th sample: $w_i = e^{\frac{-||\mathbf{x}-\mathbf{x}_i||^2}{\sigma^2}}$
 - σ = the bandwidth parameter, expresses how quickly our weight function "drops off" as points get further and further from the query x

• Classification:
$$y = \operatorname{argmax}_{c} \sum_{i=1}^{N} w_{i} I(y_{i} = c)$$

• Regression: $y = \frac{\sum_{i=1}^{N} w_{i} y_{i}}{\sum_{i=1}^{N} w_{i}}$

Another generalization: Weighted K-NNs

- Extremes
 - Bandwidth = infinity: prediction is dataset average
 - Bandwidth = zero: prediction becomes 1-NN

Kernel Regression/Classification

Four things make a memory based learner:

- A distance metric
 - Euclidean (and others)
- How many nearby neighbors to look at?
 All of them
- A weighting function (optional)
 - $w_i = exp(-d(x_i, query)^2 / \sigma^2)$
 - Nearby points to the query are weighted strongly, far points weakly.
 The σ parameter is the kernel width / bandwidth.
- How to fit with the local points?
 - Predict the weighted average of the outputs

Problems with Instance-Based Learning

- Too many features?
 - Doesn't work well if large number of irrelevant features, distances overwhelmed by noisy features
 - Distances become meaningless in high dimensions (the curse of dimensionality)
- What is the impact of the **value of K**?
- Expensive
 - No learning: most real work done during testing
 - For every test sample, must search through all dataset very slow!
 - Must use tricks like approximate nearest neighbor search
 - Need to store all training data

Curse of Dimensionality

How many neighborhoods are there?



#bins = 10x10d = 2

Atoms in the universe ~ 10⁸⁰

Subhransu Maji (UMASS)
Curse of Dimensionality

Regions become more sparsely populated given the same amount of data



Need more data to densely populate them

Figures from https://www.kdnuggets.com/2017/04/must-know-curse-dimensionality.html, https://medium.freecodecamp.org/the-curse-of-dimensionality-how-we-can-save-big-data-from-itself-d9fa0f872335

simplifies /

Increasing k <u>complicates</u> decision boundary

- Increasing k "simplifies" decision boundary
 - Majority voting means less emphasis on individual points

K = 1



K = 3



- Increasing k "simplifies" decision boundary
 - Majority voting means less emphasis on individual points

K = 5



K = 7



- Increasing k "simplifies" decision boundary
 - Majority voting means less emphasis on individual points







Summary

- K-Nearest Neighbor is the most basic and simplest to implement classifier
- Cheap at training time, expensive at test time
- Unlike other methods we'll see later, naturally works for any number of classes
- Pick K through a validation set, use approximate methods for finding neighbors
- Success of classification depends on the amount of data and the meaningfulness of the distance function (also true for other algorithms)

Plan for this lecture

- The simplest classifier: K-Nearest Neighbors
 - Algorithm and example use
 - Generalizing: Distance metrics, weighing neighbors
 - Problems: curse of dimensionality, picking K
- Logistic regression
 - Probability: review
 - Linear regression for classification?
 - Maximum likelihood solution for *logistic* regression
 - Related algorithm: perceptron

A is non-deterministic event

Can think of A as a Boolean-valued variable

Examples

- A = your next patient has cancer
- A = Steelers win Super Bowl LIII

What does P(A) mean?

Frequentist View limit N→∞ #(A is true)/N frequency of a repeating non-deterministic event

Bayesian View

P(A) is your "belief" about A

Adapted from Dhruv Batra

 $0 \le P(A) \le 1$ P(false) = 0 P(true) = 1 P(A v B) = P(A) + P(B) - P(A ^ B)



 $0 \le P(A) \le 1$ P(false) = 0 P(true) = 1 $P(A \lor B) = P(A) + P(B) - P(A \land B)$



The area of A can't get any smaller than 0

And a zero area would mean no world could ever have A true

- 0<= P(A) <= 1
- P(false) = 0
- P(true) = 1
- $P(A \lor B) = P(A) + P(B) P(A \land B)$



The area of A can'l get any bigger than 1

And an area of 1 would mean all worlds will have A true

- 0<= P(A) <= 1
- P(false) = 0
- P(true) = 1

 $P(A \lor B) = P(A) + P(B) - P(A \land B)$





Simple addition and subtraction

Dhruv Batra, Andrew Moore

Probabilities: Example Use

Apples and Oranges



Chris Bishop

Marginal, Joint, Conditional



Marginal Probability

$$p(X = x_i) = \frac{c_i}{N}.$$

Joint Probability

$$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N}$$

Conditional Probability

$$p(Y = y_j | X = x_i) = \frac{n_{ij}}{c_i}$$

Joint Probability

P(X₁,...,X_n) gives the probability of every combination of values (an *n*-dimensional array with *vⁿ* values if all variables are discrete with *v* values, all *vⁿ* values must sum to 1):

p c c c c			
	circle	square	
red	0.20	0.02	
blue	0.02	0.01	

nositive

negative

	circle	square
red	0.05	0.30
blue	0.20	0.20

 The probability of all possible conjunctions (assignments of values to some subset of variables) can be calculated by summing the appropriate subset of values from the joint distribution.

 $\begin{array}{c} P(red \land circle) \\ P(red) \end{array}$

• Therefore, all conditional probabilities can also be calculated.

 $P(positive | red \land circle)$

Adapted from Ray Mooney

Marginal Probability



P(Y=y | X=x): What do you believe about Y=y, if I tell you X=x?

P(Andy Murray wins Australian Open 2019)?

What if I tell you:

He has won it five times before He is currently ranked #307

Conditional Probability



Conditional Probability



Sum and Product Rules



Sum Rule $p(X = x_i) = \frac{c_i}{N} = \frac{1}{N} \sum_{j=1}^{L} n_{ij}$ $= \sum_{i=1}^{L} p(X = x_i, Y = y_j)$

Product Rule

$$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N} = \frac{n_{ij}}{c_i} \cdot \frac{c_i}{N}$$
$$= p(Y = y_j | X = x_i) p(X = x_i)$$

Chain Rule

Generalizes the product rule:

$$P\left(\bigcap_{k=1}^{n} A_{k}\right) = \prod_{k=1}^{n} P\left(A_{k} \middle| \bigcap_{j=1}^{k-1} A_{j}\right)$$

Example:

 $P(A_4, A_3, A_2, A_1) = P(A_4 \mid A_3, A_2, A_1) \cdot P(A_3 \mid A_2, A_1) \cdot P(A_2 \mid A_1) \cdot P(A_1)$

Equations from Wikipedia

A and B are *independent* iff:

 $P(A \mid B) = P(A)$ These two constraints are logically equivalent $P(B \mid A) = P(B)$

Therefore, if A and B are independent:

$$P(A \mid B) = \frac{P(A \land B)}{P(B)} = P(A)$$

 $P(A \land B) = P(A)P(B)$

Marginal: P satisfies $(X \perp Y)$ if and only if P(X=x,Y=y) = P(X=x) P(Y=y), $\forall x \in Val(X), y \in Val(Y)$

Conditional: *P* satisfies $(X \perp Y \mid Z)$ if and only if P(X,Y|Z) = P(X|Z) P(Y|Z), $\forall x \in Val(X), y \in Val(Y), z \in Val(Z)$

Independence

P(x,y)





Bayes' Theorem

$$p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)}$$

$$p(X) = \sum_{Y} p(X|Y)p(Y)$$

posterior ∞ likelihood × prior

Expectations

$$\mathbb{E}[f] = \sum_{x} p(x) f(x)$$

$$\mathbb{E}[f] = \int p(x)f(x) \,\mathrm{d}x$$

$$\mathbb{E}_{x}[f|y] = \sum_{x} p(x|y)f(x)$$

Conditional Expectation (discrete)

$$\mathbb{E}[f] \simeq \frac{1}{N} \sum_{n=1}^{N} f(x_n)$$

Approximate Expectation (discrete and continuous)

Entropy

$$\mathbf{H}[x] = -\sum_{x} p(x) \log_2 p(x)$$

Important quantity in

- coding theory
- statistical physics
- machine learning

Entropy



The Kullback-Leibler Divergence

$$\begin{aligned} \mathrm{KL}(p \| q) &= -\int p(\mathbf{x}) \ln q(\mathbf{x}) \, \mathrm{d}\mathbf{x} - \left(-\int p(\mathbf{x}) \ln p(\mathbf{x}) \, \mathrm{d}\mathbf{x} \right) \\ &= -\int p(\mathbf{x}) \ln \left\{ \frac{q(\mathbf{x})}{p(\mathbf{x})} \right\} \, \mathrm{d}\mathbf{x} \end{aligned}$$

$$\mathrm{KL}(p\|q) \simeq \frac{1}{N} \sum_{n=1}^{N} \left\{ -\ln q(\mathbf{x}_n | \boldsymbol{\theta}) + \ln p(\mathbf{x}_n) \right\}$$

 $\operatorname{KL}(p||q) \ge 0$ $\operatorname{KL}(p||q) \not\equiv \operatorname{KL}(q||p)$

Chris Bishop

Mutual Information

$$I[\mathbf{x}, \mathbf{y}] \equiv KL(p(\mathbf{x}, \mathbf{y}) || p(\mathbf{x}) p(\mathbf{y}))$$

= $-\iint p(\mathbf{x}, \mathbf{y}) \ln \left(\frac{p(\mathbf{x}) p(\mathbf{y})}{p(\mathbf{x}, \mathbf{y})} \right) d\mathbf{x} d\mathbf{y}$

$$I[\mathbf{x}, \mathbf{y}] = H[\mathbf{x}] - H[\mathbf{x}|\mathbf{y}] = H[\mathbf{y}] - H[\mathbf{y}|\mathbf{x}]$$

Likelihood / Prior / Posterior

- A hypothesis (model, function, parameter set, weights) is denoted as *h*; it is one member of the hypothesis space *H*
- A set of training examples is denoted as *D*, a collection of (**x**, y) pairs for training
- Pr(h) the prior probability of the hypothesis without observing any training data, what is the probability that h is the target function we want?

Likelihood / Prior / Posterior

- Pr(D) the prior probability of the observed data
 chance of getting the particular set of training examples D
- Pr(h|D) the posterior probability of h what is the probability that h is the target given that we have observed D?
- Pr(D|h) the probability of getting D if h were true (a.k.a. *likelihood of the data*)
- Pr(h|D) = Pr(D|h)Pr(h)/Pr(D)

Maximum likelihood estimation (MLE): $h_{MI} = \operatorname{argmax} Pr(D|h)$

Maximum-a-posteriori (MAP) estimation: h_{MAP} = argmax_h Pr(h|D) = argmax_h Pr(D|h)Pr(h)/Pr(D) = argmax_h Pr(D|h)Pr(h)

Classification via regression

- Suppose we ignore the fact that the target output y is binary (e.g., 0/1) rather than a continuous variable
- So we will estimate a linear regression function

$$f(\mathbf{x}; \mathbf{w}) = w_0 + w_1 x_1 + \ldots + w_d x_d$$
$$= w_0 + \mathbf{x}^T \mathbf{w}_1,$$

based on the available data as before.

• Objective we want to minimize:

$$J_n(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i; \mathbf{w}))^2$$
• We can use the resulting regression function

$$f(\mathbf{x}; \hat{\mathbf{w}}) = w_0 + \mathbf{x}^T \hat{\mathbf{w}}_1,$$

to classify any new (test) example \mathbf{x} according to

label = 1 if $f(\mathbf{x}; \mathbf{w}) > 0.5$, and label = 0 otherwise

f(x; ŵ) = 0.5 therefore defines a linear decision boundary that partitions the input space into two class specific regions (half spaces)

 Given the dissociation between the objective (classification) and the estimation criterion (regression) it is not clear that this approach leads to sensible results







Figures adapted from from Andrew Ng

The effect of outliers: Another example



Figures from Bishop

Logistic regression

- Also has "regression" in the name but it's a method for classification
- Also uses a linear combination of the features to predict label, but in a slightly different way
- Fit a *sigmoid* function to model the *probability* of the data belonging to a certain class



Background: simple decision theory

• Suppose we know the class-conditional densities $p(\mathbf{x}|y)$ for y = 0, 1 as well as the overall class frequencies P(y).

How do we decide which class a new example \mathbf{x}' belongs to so as to minimize the overall probability of error?



Background: simple decision theory

• Suppose we know the class-conditional densities $p(\mathbf{x}|y)$ for y = 0, 1 as well as the overall class frequencies P(y).

How do we decide which class a new example \mathbf{x}' belongs to so as to minimize the overall probability of error?



The minimum probability of error decisions are given by

$$y' = \arg \max_{y=0,1} \{ p(\mathbf{x}'|y)P(y) \}$$
$$= \arg \max_{y=0,1} \{ P(y|\mathbf{x}') \}$$

Logistic regression

• The optimal decisions are based on the posterior class probabilities $P(y|\mathbf{x})$. For binary classification problems, we can write these decisions as

$$y = 1$$
 if $\log \frac{P(y = 1 | \mathbf{x})}{P(y = 0 | \mathbf{x})} > 0$

and y = 0 otherwise.

Logistic regression

• The optimal decisions are based on the posterior class probabilities $P(y|\mathbf{x})$. For binary classification problems, we can write these decisions as

$$y = 1$$
 if $\log \frac{P(y = 1 | \mathbf{x})}{P(y = 0 | \mathbf{x})} > 0$

and y = 0 otherwise.

• We generally don't know $P(y|\mathbf{x})$ but we can parameterize the possible decisions according to

$$\log \frac{P(y=1|\mathbf{x})}{P(y=0|\mathbf{x})} = f(\mathbf{x};\mathbf{w}) = w_0 + \mathbf{x}^T \mathbf{w}_1$$

Logistic regression cont'd

• Our log-odds model

$$\log \frac{P(y=1|\mathbf{x})}{P(y=0|\mathbf{x})} = w_0 + \mathbf{x}^T \mathbf{w}_1$$

gives rise to a specific form for the conditional probability over the labels (the logistic model):

$$P(y = 1 | \mathbf{x}, \mathbf{w}) = \boldsymbol{\sigma} (w_0 + \mathbf{x}^T \mathbf{w}_1)$$

where

$$\sigma(z) = (1 + \exp(-z))^{-1}$$

is a logistic "squashing function" that turns linear predictions into probabilities



Logistic regression: decisions

• Logistic regression models imply a linear decision boundary

$$\log \frac{P(y=1|\mathbf{x})}{P(y=0|\mathbf{x})} = w_0 + \mathbf{x}^T \mathbf{w}_1 = 0$$



Fitting logistic regression models

 We can fit the logistic models using the maximum (conditional) log-likelihood criterion

$$l(D; \mathbf{w}) = \sum_{i=1}^{n} \log P(y_i | \mathbf{x}_i, \mathbf{w})$$

where

$$P(y = 1 | \mathbf{x}, \mathbf{w}) = \boldsymbol{\sigma} \left(w_0 + \mathbf{x}^T \mathbf{w}_1 \right)$$

• Solution: find roots of
$$(y_i - P(y_i = 1 | \mathbf{x}_i, \mathbf{w})) \begin{bmatrix} 1 \\ \mathbf{x}_i \end{bmatrix} = 0$$
 prediction error

Stochastic gradient ascent

 We can try to maximize the log-likelihood in an *on-line* or incremental fashion.

Given each training input \mathbf{x}_i and the binary (0/1) label y_i , we can change the parameters \mathbf{w} slightly to increase the corresponding log-probability

$$\mathbf{w} \leftarrow \mathbf{w} + \eta \frac{\partial}{\partial \mathbf{w}} \log P(y_i | \mathbf{x}_i, \mathbf{w})$$

= $\mathbf{w} + \eta \underbrace{\left(y_i - P(y_i = 1 | \mathbf{x}_i, \mathbf{w})\right)}_{\text{prediction error}} \begin{bmatrix} 1 \\ \mathbf{x}_i \end{bmatrix}$

where η is the *learning rate*.

Whiteboard: solution

Logistic Regression / MLE Example

- Want to find the weight vector that gives us the highest P(y_i|x_i, w),
- where P(y_i=1|x_i, w) = 1 / (1 + exp(-w'*x))
- Consider two weight vectors and three samples, with corresponding likelihoods:

	P(y_1 = 1 x_1, w_i)	P(y_2 = 1 x_2, w_i)	P(y_3 = 1 x_3, w_i)
w_ 1	0.3	0.1	0.4
w_ 2	0.7	0.8	0.2
True label:	1	0	1

Logistic Regression / MLE Example

- Then the value of the objective for w_i is: P(y_1 = 1 | x_1, w_i) * (1 - P(y_2 = 1 | x_2, w_i)) * P(y_3 = 1 | x_3, w_i)
- So the score for w_1 is: 0.3 * 0.9 * 0.4
- And the score for w_2 is: 0.7 * 0.2 * 0.2
- Thus, w_1 is a better weight vector = model

Plan for this lecture

- The simplest classifier: K-Nearest Neighbors
 - Algorithm and example use
 - Generalizing: Distance metrics, weighing neighbors
 - Problems: curse of dimensionality, picking K
- Logistic regression
 - Probability: review
 - Linear regression for classification?
 - Maximum likelihood solution for *logistic* regression
 - Related algorithm: perceptron

The perceptron algorithm

- Rosenblatt (1962)
- Prediction rule: $y(\mathbf{x}) = f(\mathbf{w}^{\mathrm{T}} \phi(\mathbf{x}))$ $f(a) = \begin{cases} +1, & a \ge 0\\ -1, & a < 0 \end{cases}$ where
- Want: $\mathbf{w}^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}_n) t_n > 0$ ($t_n = +1 \text{ or } -1$) Loss: $E_{\mathrm{P}}(\mathbf{w}) = -\sum \mathbf{w}^{\mathrm{T}} \boldsymbol{\phi}_n t_n$
- $n \in \mathcal{M}$

(just using the Misclassified examples)

The perceptron algorithm

• Loss:
$$E_{\mathrm{P}}(\mathbf{w}) = -\sum_{n \in \mathcal{M}} \mathbf{w}^{\mathrm{T}} \boldsymbol{\phi}_n t_n$$

• Learning algorithm update rule:

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_{\mathrm{P}}(\mathbf{w}) = \mathbf{w}^{(\tau)} + \eta \phi_n t_n$$

- Interpretation:
 - If sample is misclassified and is positive, make the weight vector more like it
 - If sample is misclassified and negative... unlike it

The perceptron algorithm (red=pos)



Figures from Bishop

Summary: Tradeoffs of classification methods thus far

- Nearest neighbors
 - Non-parametric method; basic formulation cannot ignore/focus on different feature dimensions
 - Slow at test time (large search problem to find neighbors)
 - Need to store all data points (unlike SVM, coming next)
 - Decision boundary not necessarily linear
 - Naturally handles multiple classes
- Logistic regression (a *classification* method)
 - Models the *probability* of a label given the data
 - Decision boundary corresponds to $\mathbf{w}^T \mathbf{x} = 0$ (a line)
- Perceptron
 - Same decision boundary as logistic regression (a line)
 - Simple update rule
 - Won't converge for non-linearly-separable data