

Linear Classifiers

André Martins



Unbabel



instituto de
telecomunicações



TÉCNICO
LISBOA

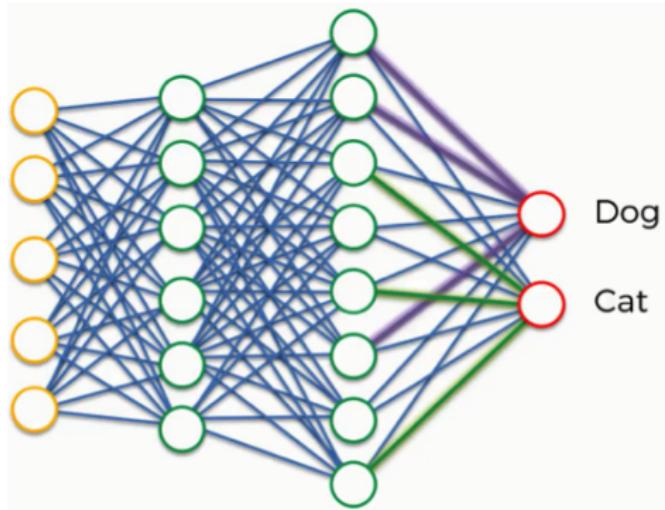
Lisbon Machine Learning School, July 12, 2019

Why Linear Classifiers?

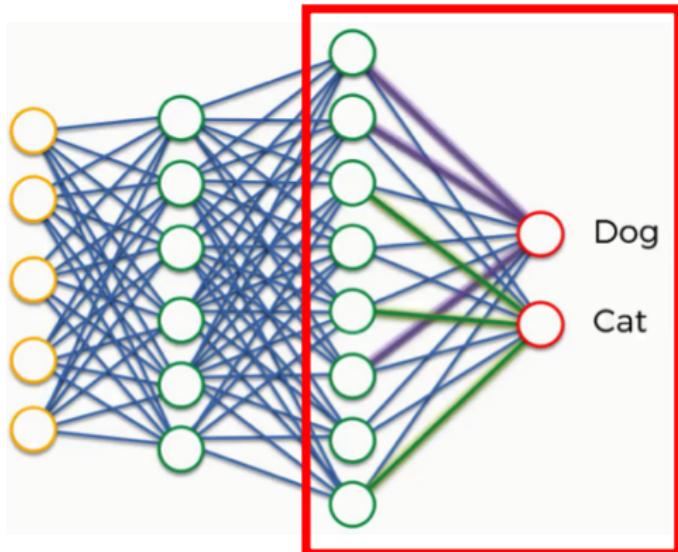
It's 2019 and everybody uses neural networks. Why a lecture on linear classifiers?

- The underlying machine learning concepts are the same
- The theory (statistics and optimization) are much better understood
- Linear classifiers are still widely used (and very effective when data is scarce)
- Linear classifiers are **a component of neural networks.**

Linear Classifiers and Neural Networks

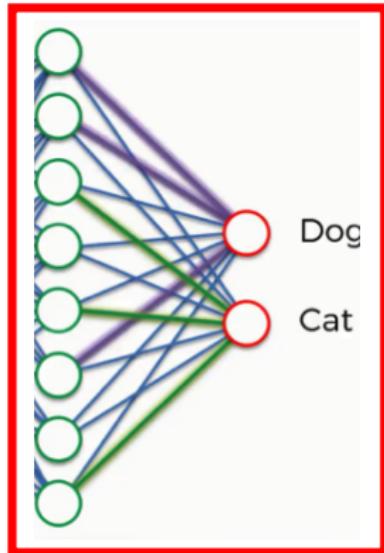


Linear Classifiers and Neural Networks



Linear Classifier

Linear Classifiers and Neural Networks

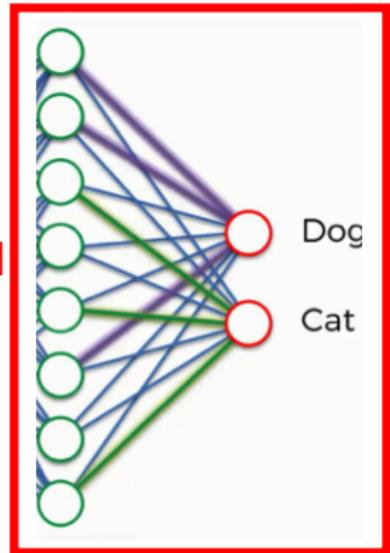


Linear Classifier

Linear Classifiers and Neural Networks



**Handcrafted
Features**



Linear Classifier

Today's Roadmap

- Binary and multi-class classification
- Linear classifiers: perceptron, naive Bayes, logistic regression, SVMs
- Softmax and sparsemax
- Regularization and optimization, stochastic gradient descent
- Similarity-based classifiers and kernels.

Fake News Detection

Task: tell if a news article / quote is **fake** or **real**.

This is a **binary classification problem**.

Fake Or Real?



Fake Or Real?

*With Artificial
Intelligence we
are summoning
the demons*
- Elon Musk



Fake Or Real?



AlphaGo Beats Go Human Champ: Godfather Of Deep Learning Tells Us Do Not Be Afraid Of AI

21 March 2016, 10:16 am EDT By [Aaron Mamiit](#) Tech Times



Last week, Google's artificial intelligence program

Last week, Google's artificial intelligence program AlphaGo [dominated](#) its match with South Korean world Go champion Lee Sedol, winning with a 4-1 score.

The achievement stunned artificial intelligence experts, who previously thought that Google's computer program would need at least 10 more years before developing enough to be able to beat a human world champion.

Fake Or Real?

The image shows the Google Moon interface. At the top left is the Google logo with "Moon" underneath. Below it is the title "Moon". On the left, there is a zoom control with a map icon, a plus sign, and a minus sign. The main area displays a yellowish-brown surface covered in numerous craters, representing the Moon's terrain. A white callout box with a black border and rounded corners is positioned in the upper-left quadrant of the map. It contains the text "Apollo 11" at the top, followed by "Jul 20, 1969" and a list of crew members: "Neil A. Armstrong, Commander; Edwin E. Aldrin, Lunar Module Pilot; Michael Collins, Command Module Pilot". A small red location pin is placed on the map near the landing site. To the right of the map, there is a sidebar with the heading "Welcome to Google Moon". Below it, a paragraph of text reads: "In honor of the first manned Moon landing, which took place on July 20, 1969, we've added some NASA imagery to the [Google Maps](#) interface to help you pay your own visit to our celestial neighbor. Happy lunar surfing. [More about Google Moon](#)". At the bottom right of the map, there is a link "Looking for something on [Planet Earth](#)?". On the far right of the sidebar, there is a vertical list of six Apollo missions, each with a red location pin icon and the mission name, date, and year:

- Apollo 11 Jul 20, 1969
- Apollo 12 Nov 19, 1969
- Apollo 14 Feb 5, 1971
- Apollo 15 Jul 30, 1971
- Apollo 16 Apr 20, 1972
- Apollo 17 Dec 11, 1972

©2009 Google - Images ©2009 NASA - Terms of Use

Fake Or Real?

Can a machine determine this automatically?

It can be a very hard problem, since fact-checking is hard and requires combining several knowledge sources

... also, reality surpasses fiction sometimes.

Topic Classification

Task: given a news article, determine its topic (politics, sports, etc.)

This is a **multi-class classification problem.**

It's a much easier task, we can get 80-90% accuracies with a simple ML model.

Topic Classification

AlphaGo Beats Go Human Champ: Godfather Of Deep Learning Tells Us Do Not Be Afraid Of AI

21 March 2016, 10:16 am EDT By Aaron Marnit Tech Times



Last week, Google's artificial intelligence program

Last week, Google's artificial intelligence program AlphaGo **dominated** its match with South Korean world Go champion Lee Sedol, winning with a 4-1 score.

The achievement stunned artificial intelligence experts, who previously thought that Google's computer program would need at least 10 more years before developing enough to be able to beat a human world champion.



sports
politics
technology
economy
weather
culture

Outline

① Data and Feature Representation

② Perceptron

③ Naive Bayes

④ Logistic Regression

⑤ Support Vector Machines

⑥ Regularization

⑦ Non-Linear Classifiers

Disclaimer

Many of the following slides are adapted from Ryan McDonald.

Let's Start Simple

- Example 1 – sequence: $\star \diamond \circ;$ label: -1
- Example 2 – sequence: $\star \heartsuit \triangle;$ label: -1
- Example 3 – sequence: $\star \triangle \spadesuit;$ label: $+1$
- Example 4 – sequence: $\diamond \triangle \circ;$ label: $+1$

Let's Start Simple

- Example 1 – sequence: $\star \diamond \circ$; label: -1
- Example 2 – sequence: $\star \heartsuit \triangle$; label: -1
- Example 3 – sequence: $\star \triangle \spadesuit$; label: +1
- Example 4 – sequence: $\diamond \triangle \circ$; label: +1

- New sequence: $\star \diamond \circ$; label ?

Let's Start Simple

- Example 1 – sequence: $\star \diamond \circ$; label: -1
- Example 2 – sequence: $\star \heartsuit \triangle$; label: -1
- Example 3 – sequence: $\star \triangle \spadesuit$; label: $+1$
- Example 4 – sequence: $\diamond \triangle \circ$; label: $+1$

- New sequence: $\star \diamond \circ$; label -1
- New sequence: $\star \diamond \heartsuit$; label ?

Let's Start Simple

- Example 1 – sequence: $\star \diamond \circ$; label: -1
 - Example 2 – sequence: $\star \heartsuit \triangle$; label: -1
 - Example 3 – sequence: $\star \triangle \spadesuit$; label: $+1$
 - Example 4 – sequence: $\diamond \triangle \circ$; label: $+1$
-
- New sequence: $\star \diamond \circ$; label -1
 - New sequence: $\star \diamond \heartsuit$; label -1
 - New sequence: $\star \triangle \circ$; label ?

Let's Start Simple

- Example 1 – sequence: $\star \diamond \circ$; label: -1
- Example 2 – sequence: $\star \heartsuit \triangle$; label: -1
- Example 3 – sequence: $\star \triangle \spadesuit$; label: $+1$
- Example 4 – sequence: $\diamond \triangle \circ$; label: $+1$

- New sequence: $\star \diamond \circ$; label -1
- New sequence: $\star \diamond \heartsuit$; label -1
- New sequence: $\star \triangle \circ$; label ?

Why can we do this?

Let's Start Simple: Machine Learning

- Example 1 – sequence: $\star \diamond \circ$; label: -1
- Example 2 – sequence: $\star \heartsuit \triangle$; label: -1
- Example 3 – sequence: $\star \triangle \spadesuit$; label: $+1$
- Example 4 – sequence: $\diamond \triangle \circ$; label: $+1$
- New sequence: $\star \diamond \heartsuit$; label -1

Label -1

Label $+1$

$$P(-1|\star) = \frac{\text{count}(\star \text{ and } -1)}{\text{count}(\star)} = \frac{2}{3} = 0.67 \text{ vs. } P(+1|\star) = \frac{\text{count}(\star \text{ and } +1)}{\text{count}(\star)} = \frac{1}{3} = 0.33$$

$$P(-1|\diamond) = \frac{\text{count}(\diamond \text{ and } -1)}{\text{count}(\diamond)} = \frac{1}{2} = 0.5 \text{ vs. } P(+1|\diamond) = \frac{\text{count}(\diamond \text{ and } +1)}{\text{count}(\diamond)} = \frac{1}{2} = 0.5$$

$$P(-1|\heartsuit) = \frac{\text{count}(\heartsuit \text{ and } -1)}{\text{count}(\heartsuit)} = \frac{1}{1} = 1.0 \text{ vs. } P(+1|\heartsuit) = \frac{\text{count}(\heartsuit \text{ and } +1)}{\text{count}(\heartsuit)} = \frac{0}{1} = 0.0$$

Let's Start Simple: Machine Learning

- Example 1 – sequence: $\star \diamond \circ$; label: -1
- Example 2 – sequence: $\star \heartsuit \triangle$; label: -1
- Example 3 – sequence: $\star \triangle \spadesuit$; label: $+1$
- Example 4 – sequence: $\diamond \triangle \circ$; label: $+1$
- New sequence: $\star \triangle \circ$; label ?

Label -1

Label $+1$

$$P(-1|\star) = \frac{\text{count}(\star \text{ and } -1)}{\text{count}(\star)} = \frac{2}{3} = 0.67 \text{ vs. } P(+1|\star) = \frac{\text{count}(\star \text{ and } +1)}{\text{count}(\star)} = \frac{1}{3} = 0.33$$

$$P(-1|\triangle) = \frac{\text{count}(\triangle \text{ and } -1)}{\text{count}(\triangle)} = \frac{1}{3} = 0.33 \text{ vs. } P(+1|\triangle) = \frac{\text{count}(\triangle \text{ and } +1)}{\text{count}(\triangle)} = \frac{2}{3} = 0.67$$

$$P(-1|\circ) = \frac{\text{count}(\circ \text{ and } -1)}{\text{count}(\circ)} = \frac{1}{2} = 0.5 \text{ vs. } P(+1|\circ) = \frac{\text{count}(\circ \text{ and } +1)}{\text{count}(\circ)} = \frac{1}{2} = 0.5$$

Machine Learning

- ① Define a model/distribution of interest
- ② Make some assumptions if needed
- ③ Fit the model to the data

Machine Learning

- ① Define a model/distribution of interest
 - ② Make some assumptions if needed
 - ③ Fit the model to the data
- Model: $P(\text{label}|\text{sequence}) = P(\text{label}|\text{symbol}_1, \dots, \text{symbol}_n)$
 - Prediction for new sequence = $\text{argmax}_{\text{label}} P(\text{label}|\text{sequence})$
 - Assumption (**naive Bayes**—more later):
$$P(\text{symbol}_1, \dots, \text{symbol}_n | \text{label}) = \prod_{i=1}^n P(\text{symbol}_i | \text{label})$$
 - Fit the model to the data: count!! (simple probabilistic modeling)

Some Notation: Inputs and Outputs

- Input $x \in \mathcal{X}$
 - e.g., a news article, a sentence, an image, ...
- Output $y \in \mathcal{Y}$
 - e.g., fake/not fake, a topic, a parse tree, an image segmentation
- Input/Output pair: $(x, y) \in \mathcal{X} \times \mathcal{Y}$
 - e.g., a **news article** together with a **topic**
 - e.g., a **sentence** together with a **parse tree**
 - e.g., an **image** partitioned into **segmentation regions**

Supervised Machine Learning

- We are given a **labeled dataset** of input/output pairs:

$$\mathcal{D} = \{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1}^N \subseteq \mathcal{X} \times \mathcal{Y}$$

- **Goal:** use it to learn a **classifier** $h : \mathcal{X} \rightarrow \mathcal{Y}$ that generalizes well to arbitrary inputs.
- At test time, given $\mathbf{x} \in \mathcal{X}$, we predict

$$\hat{\mathbf{y}} = h(\mathbf{x}).$$

- Hopefully, $\hat{\mathbf{y}} \approx \mathbf{y}$ most of the time.

Things can go by different names depending on what \mathcal{Y} is...

Regression

Deals with **continuous** output variables:

- **Regression:** $\mathcal{Y} = \mathbb{R}$
 - e.g., given a news article, how much time a user will spend reading it?
- **Multivariate regression:** $\mathcal{Y} = \mathbb{R}^K$
 - e.g., predict the X-Y coordinates in an image where the user will click

Classification

Deals with **discrete** output variables:

- **Binary classification:** $\mathcal{Y} = \{\pm 1\}$
 - e.g., fake news detection
- **Multi-class classification:** $\mathcal{Y} = \{1, 2, \dots, K\}$
 - e.g., topic classification
- **Structured classification:** \mathcal{Y} exponentially large and structured
 - e.g., machine translation, caption generation, image segmentation

At LxMLS: we'll get into **structured classification**

... but for now, let's talk about multi-class classification first.

Sometimes **reductions** are convenient:

- logistic regression reduces classification to regression
- one-vs-all reduces multi-class to binary
- greedy search reduces structured classification to multi-class

... but other times it's better to tackle the problem in its native form.

More later!

Feature Representations

Feature engineering is an important step in linear classifiers:

- Bag-of-words features for text, also lemmas, parts-of-speech, ...
- SIFT features and wavelet representations in computer vision
- Other categorical, Boolean, and continuous features

Feature Representations

We need to represent information about x

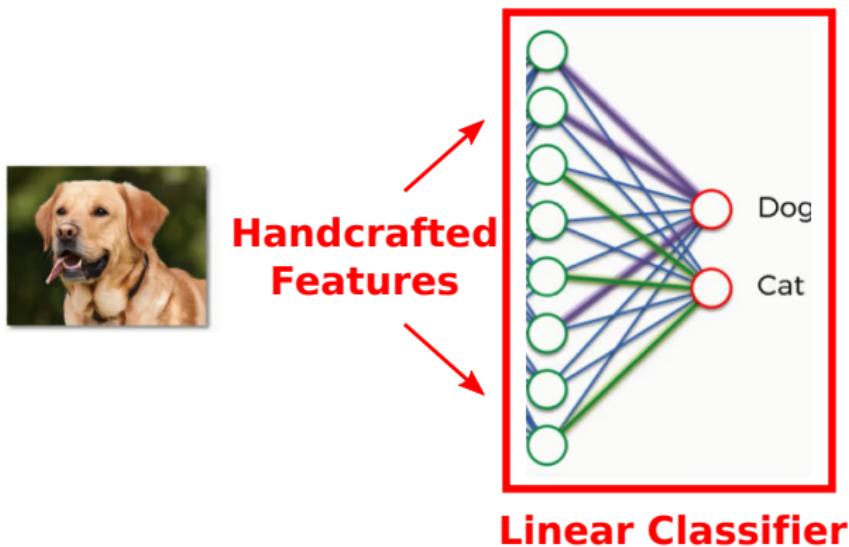
Typical approach: define a feature map $\psi : \mathcal{X} \rightarrow \mathbb{R}^D$

- $\psi(x)$ is a high dimensional **feature vector**

We can use feature vectors to encapsulate **Boolean**, **categorical**, and **continuous** features

- e.g., categorical features can be reduced to a range of one-hot binary values.

Example: Continuous Features



Feature Representations: Joint Feature Mappings

For multi-class/structured classification, a **joint feature map**
 $\phi : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^D$ is sometimes more convenient

- $\phi(x, y)$ instead of $\psi(x)$

Each feature now represents a joint property of the input x and the candidate output y .

We'll use this notation in the labs this afternoon!

Examples

- x is a document and y is a label

$$\phi_j(x, y) = \begin{cases} 1 & \text{if } x \text{ contains the word "interest"} \\ & \text{and } y = \text{"financial"} \\ 0 & \text{otherwise} \end{cases}$$

$\phi_j(x, y) = \%$ of words in x with punctuation and $y = \text{"scientific"}$

- x is a word and y is a part-of-speech tag

$$\phi_j(x, y) = \begin{cases} 1 & \text{if } x = \text{"bank" and } y = \text{Verb} \\ 0 & \text{otherwise} \end{cases}$$

More Examples

- x is a name, y is a label classifying the type of entity

$$\phi_0(x, y) = \begin{cases} 1 & \text{if } x \text{ contains "George"} \\ & \text{and } y = \text{"Person"} \\ 0 & \text{otherwise} \end{cases}$$

$$\phi_4(x, y) = \begin{cases} 1 & \text{if } x \text{ contains "George"} \\ & \text{and } y = \text{"Location"} \\ 0 & \text{otherwise} \end{cases}$$

$$\phi_1(x, y) = \begin{cases} 1 & \text{if } x \text{ contains "Washington"} \\ & \text{and } y = \text{"Person"} \\ 0 & \text{otherwise} \end{cases}$$

$$\phi_5(x, y) = \begin{cases} 1 & \text{if } x \text{ contains "Washington"} \\ & \text{and } y = \text{"Location"} \\ 0 & \text{otherwise} \end{cases}$$

$$\phi_2(x, y) = \begin{cases} 1 & \text{if } x \text{ contains "Bridge"} \\ & \text{and } y = \text{"Person"} \\ 0 & \text{otherwise} \end{cases}$$

$$\phi_6(x, y) = \begin{cases} 1 & \text{if } x \text{ contains "Bridge"} \\ & \text{and } y = \text{"Location"} \\ 0 & \text{otherwise} \end{cases}$$

$$\phi_3(x, y) = \begin{cases} 1 & \text{if } x \text{ contains "General"} \\ & \text{and } y = \text{"Person"} \\ 0 & \text{otherwise} \end{cases}$$

$$\phi_7(x, y) = \begin{cases} 1 & \text{if } x \text{ contains "General"} \\ & \text{and } y = \text{"Location"} \\ 0 & \text{otherwise} \end{cases}$$

- $x=\text{General George Washington}$, $y=\text{Person}$ $\rightarrow \phi(x, y) = [1 \ 1 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0]$
- $x=\text{George Washington Bridge}$, $y=\text{Location}$ $\rightarrow \phi(x, y) = [0 \ 0 \ 0 \ 0 \ 1 \ 1 \ 1 \ 0]$
- $x=\text{George Washington George}$, $y=\text{Location}$ $\rightarrow \phi(x, y) = [0 \ 0 \ 0 \ 0 \ 1 \ 1 \ 0 \ 0]$

Block Feature Vectors

- $x = \text{General George Washington}$, $y = \text{Person} \rightarrow \phi(x, y) = [1 \ 1 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0]$
- $x = \text{General George Washington}$, $y = \text{Location} \rightarrow \phi(x, y) = [0 \ 0 \ 0 \ 0 \ 1 \ 1 \ 0 \ 1]$
- $x = \text{George Washington Bridge}$, $y = \text{Location} \rightarrow \phi(x, y) = [0 \ 0 \ 0 \ 0 \ 1 \ 1 \ 1 \ 0]$
- $x = \text{George Washington George}$, $y = \text{Location} \rightarrow \phi(x, y) = [0 \ 0 \ 0 \ 0 \ 1 \ 1 \ 0 \ 0]$
- Each equal size block of the feature vector corresponds to one label
- Non-zero values allowed only in one block

Feature Representations – $\psi(x)$ vs. $\phi(x, y)$

Equivalent if $\phi(x, y)$ conjoins input features $\psi(x)$ with one-hot label representations $e_y := [0, \dots, 0, 1, 0, \dots, 0]$

$$\begin{aligned}\phi(x, y) &= \psi(x) \otimes e_y \\ &= [0, \dots, 0, \underbrace{\psi(x)}_{y^{\text{th}} \text{ block}}, 0, \dots, 0]\end{aligned}$$

- $\psi(x)$
 - $x = \text{General George Washington} \rightarrow \psi(x) = [1 \ 1 \ 0 \ 1]$
- $\phi(x, y)$
 - $x = \text{General George Washington}, y = \text{Person} \rightarrow \phi(x, y) = [1 \ 1 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0]$
 - $x = \text{General George Washington}, y = \text{Object} \rightarrow \phi(x, y) = [0 \ 0 \ 0 \ 0 \ 1 \ 1 \ 0 \ 1]$

$\psi(x)$ is sometimes simpler and more convenient in binary classification
... but $\phi(x, y)$ is more expressive (allows more complex features over properties of labels)

Feature Engineering and NLP Pipelines

Classical NLP pipelines consist of stacking together several linear classifiers

Each classifier's predictions are used to handcraft features for other classifiers

Examples of features:

- POS tags: adjective counts for sentiment analysis
- Spell checker: misspellings counts for spam detection
- Parsing: depth of tree for readability assessment.

Example: Translation Quality Estimation

The screenshot shows the Google Translate interface. At the top, there's a navigation bar with the Google logo, a grid icon, a bell icon, and a profile picture. Below it, the word "Translate" is written in red, with a "Turn off instant translation" link and a star icon. The main area has two language selection dropdowns: one for the source language (English, Spanish, French, Detect language) and one for the target language (French, Spanish, Portuguese). A "Translate" button is located between them. On the left, the input text "does machine translation work?" is shown in English, with a character count of 30/5000. On the right, the translated text "Le travail de traduction automatique?" is shown in French. Both text boxes have edit icons at the bottom. At the very bottom, there's a navigation bar with icons for back, forward, search, and other functions.

Google

Translate Turn off instant translation

English Spanish French Detect language

French Spanish Portuguese

Translate

does machine translation work? x

30/5000

Le travail de traduction automatique?

Andre Martins (IST)

Linear Classifiers

LxMLS 2019

33 / 144

Example: Translation Quality Estimation

Wrong translation!

The screenshot shows the Google Translate interface. On the left, the input text "does machine translation work?" is displayed in English. On the right, the output translation "Le travail de traduction automatique?" is shown in French. A red oval highlights the French translation, and a red arrow points from the text "Wrong translation!" at the top to this highlighted area. The interface includes language selection dropdowns (English to French), a "Translate" button, and various input/output options like microphone and keyboard icons.

Example: Translation Quality Estimation

Wrong translation!

A screenshot of the Google Translate interface. The input text "does machine translation work?" is in English. The target language is set to French, and the output is "Le travail de traduction automatique?". A red oval highlights the output text, and a red arrow points from the text "Wrong translation!" above to the highlighted area. The interface includes language selection dropdowns, a "Translate" button, and various interaction icons.

Goal: estimate the quality of a translation on the fly (without a reference)!

Example: Translation Quality Estimation

Hand-crafted features:

- no of tokens in the source/target segment
- LM probability of source/target segment and their ratio
- % of source 1–3-grams observed in 4 frequency quartiles of source corpus
- average no of translations per source word
- ratio of brackets and punctuation symbols in source & target segments
- ratio of numbers, content/non-content words in source & target segments
- ratio of nouns/verbs/etc in the source & target segments
- % of dependency relations b/w constituents in source & target segments
- diff in depth of the syntactic trees of source & target segments
- diff in no of PP/NP/VP/ADJP/ADVP/CONJP in source & target
- diff in no of person/location/organization entities in source & target
- features and global score of the SMT system
- number of distinct hypotheses in the n-best list
- 1–3-gram LM probabilities using translations in the n-best to train the LM
- average size of the target phrases
- proportion of pruned search graph nodes;
- proportion of recombined graph nodes.

Representation Learning

Feature engineering is a black art and can be very time-consuming
But it's a good way of encoding prior knowledge, and it is still widely used
in practice (in particular with “small data”)
One alternative to feature engineering: **representation learning**

Bhiksha will discuss this tomorrow!

Our Setup

Let's assume a multi-class classification problem, with $|\mathcal{Y}|$ labels (classes).

Linear Classifiers

- Parametrized by a **weight vector** $w \in \mathbb{R}^D$ (one weight per feature)
- The score (or probability) of a particular label is based on a **linear** combination of features and their weights
- At test time (known w), predict the class \hat{y} which maximizes this score:

$$\hat{y} = h(x) = \arg \max_{y \in \mathcal{Y}} w \cdot \phi(x, y)$$

- At training time, different strategies to learn w yield different linear classifiers: perceptron, naïve Bayes, logistic regression, SVMs, ...

Linear Classifiers – $\psi(x)$

- Define $|\mathcal{Y}|$ weight vectors $w_y \in \mathbb{R}^D$
 - i.e., one weight vector per output label y
- **Classification**

$$\hat{y} = \arg \max_{y \in \mathcal{Y}} w_y \cdot \psi(x)$$

Linear Classifiers – $\psi(x)$

- Define $|\mathcal{Y}|$ weight vectors $w_y \in \mathbb{R}^D$
 - i.e., one weight vector per output label y

- **Classification**

$$\hat{y} = \arg \max_{y \in \mathcal{Y}} w_y \cdot \psi(x)$$

- $\phi(x, y)$
 - $x=\text{General George Washington}, y=\text{Person} \rightarrow \phi(x, y) = [1 \ 1 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0]$
 - $x=\text{General George Washington}, y=\text{Object} \rightarrow \phi(x, y) = [0 \ 0 \ 0 \ 0 \ 1 \ 1 \ 0 \ 1]$
 - Single $w \in \mathbb{R}^8$
- $\psi(x)$
 - $x=\text{General George Washington} \rightarrow \psi(x) = [1 \ 1 \ 0 \ 1]$
 - Two parameter vectors $w_0 \in \mathbb{R}^4, w_1 \in \mathbb{R}^4$

Linear Classifiers – Bias Terms

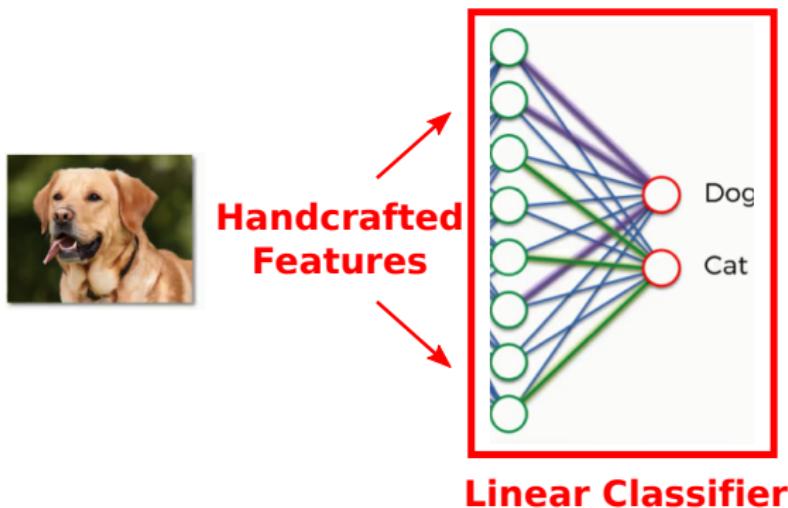
- Often linear classifiers are presented as

$$\hat{y} = \arg \max_{y \in \mathcal{Y}} w_y \cdot \psi(x) + b_y$$

where b_y is a bias or offset term

- This can be folded into $\psi(x)$ (by defining a constant feature for each label)
- We assume this for simplicity.

Commonly Used Notation in Neural Networks



$$\hat{y} = \text{argmax} (\mathbf{W}\psi(x) + \mathbf{b}), \quad \mathbf{W} = \begin{bmatrix} \vdots \\ w_y^\top \\ \vdots \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} \vdots \\ b_y \\ \vdots \end{bmatrix}.$$

Binary Linear Classifier

With **binary labels** ($\mathcal{Y} = \{\pm 1\}$) we often use a minimal parametrization:

$$\hat{y} = \arg \max_{y \in \{\pm 1\}} w_y \cdot \psi(x) + b_y$$

Binary Linear Classifier

With **binary labels** ($\mathcal{Y} = \{\pm 1\}$) we often use a minimal parametrization:

$$\begin{aligned}\hat{y} &= \arg \max_{y \in \{\pm 1\}} w_y \cdot \psi(x) + b_y \\ &= \begin{cases} +1 & \text{if } w_{+1} \cdot \psi(x) + b_{+1} > w_{-1} \cdot \psi(x) + b_{-1} \\ -1 & \text{otherwise} \end{cases}\end{aligned}$$

Binary Linear Classifier

With **binary labels** ($\mathcal{Y} = \{\pm 1\}$) we often use a minimal parametrization:

$$\begin{aligned}\hat{y} &= \arg \max_{y \in \{\pm 1\}} \mathbf{w}_y \cdot \psi(x) + b_y \\ &= \begin{cases} +1 & \text{if } \mathbf{w}_{+1} \cdot \psi(x) + b_{+1} > \mathbf{w}_{-1} \cdot \psi(x) + b_{-1} \\ -1 & \text{otherwise} \end{cases} \\ &= \text{sign}(\underbrace{(\mathbf{w}_{+1} - \mathbf{w}_{-1}) \cdot \psi(x)}_v + \underbrace{(b_{+1} - b_{-1})}_c).\end{aligned}$$

Binary Linear Classifier

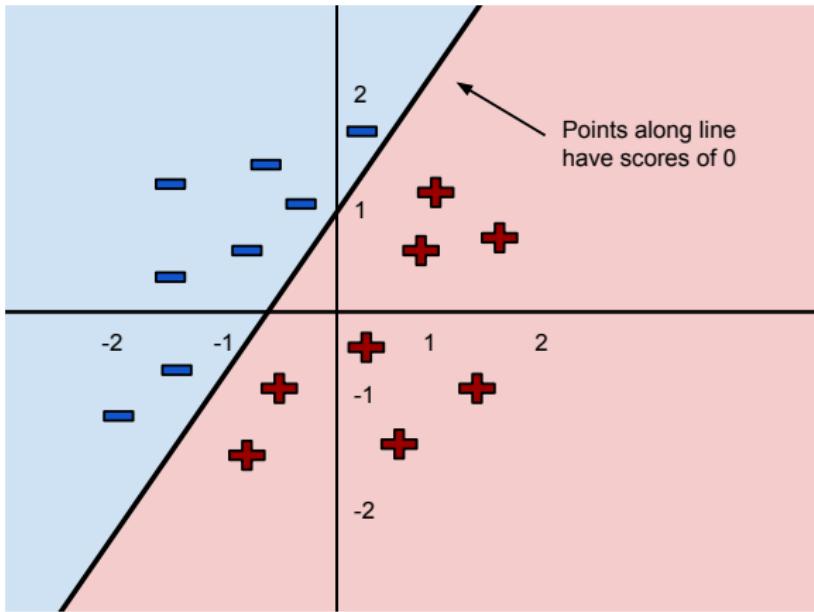
With **binary labels** ($\mathcal{Y} = \{\pm 1\}$) we often use a minimal parametrization:

$$\begin{aligned}\hat{y} &= \arg \max_{y \in \{\pm 1\}} \mathbf{w}_y \cdot \psi(x) + b_y \\ &= \begin{cases} +1 & \text{if } \mathbf{w}_{+1} \cdot \psi(x) + b_{+1} > \mathbf{w}_{-1} \cdot \psi(x) + b_{-1} \\ -1 & \text{otherwise} \end{cases} \\ &= \text{sign}(\underbrace{(\mathbf{w}_{+1} - \mathbf{w}_{-1}) \cdot \psi(x)}_v + \underbrace{(b_{+1} - b_{-1})}_c).\end{aligned}$$

That is: only half of the parameters are needed.

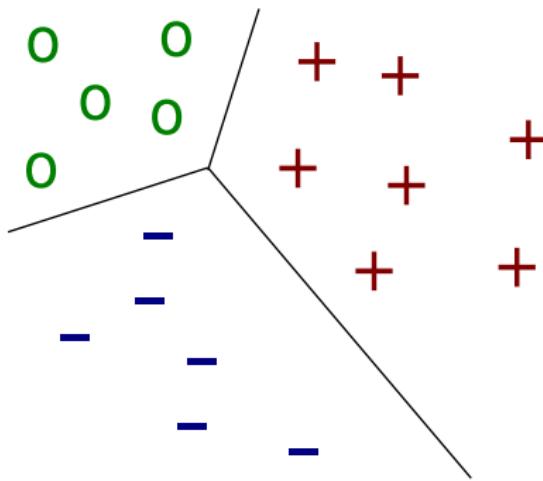
Binary Linear Classifier

Then (v, c) is an hyperplane that divides all points:



Multiclass Linear Classifier

Defines regions of space.

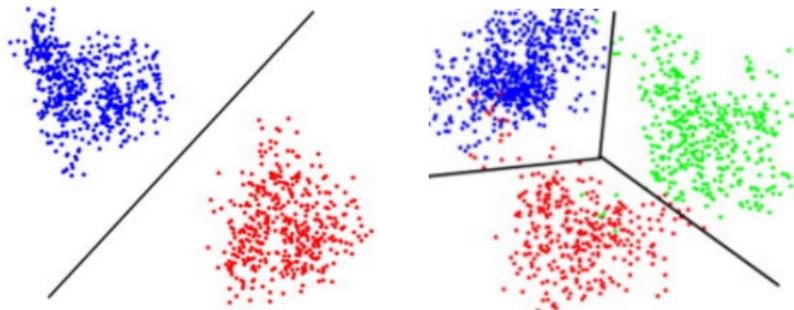


Linear Classifiers

- Prediction rule:

$$\hat{y} = h(x) = \arg \max_{y \in \mathcal{Y}} \overbrace{\mathbf{w} \cdot \phi(x, y)}^{\text{linear in } \mathbf{w}}$$

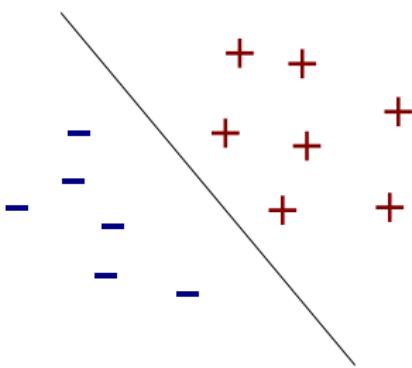
- The decision boundary is defined by the intersection of half spaces
- In the binary case ($|\mathcal{Y}| = 2$) this corresponds to a hyperplane classifier



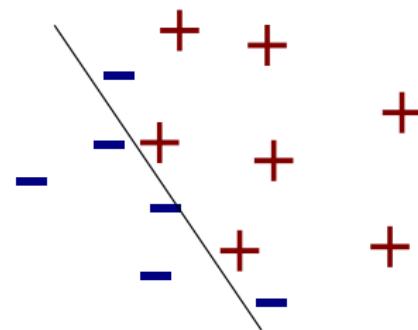
Linear Separability

- A set of points is **linearly separable** if there exists a w such that classification is perfect

Separable



Not Separable



Outline

① Data and Feature Representation

② Perceptron

③ Naive Bayes

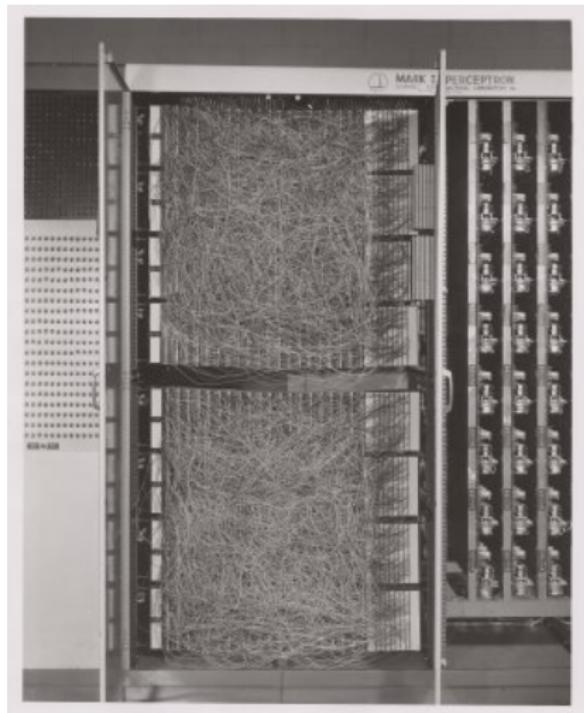
④ Logistic Regression

⑤ Support Vector Machines

⑥ Regularization

⑦ Non-Linear Classifiers

Perceptron (Rosenblatt, 1958)



(Extracted from Wikipedia)

- Invented in 1957 at the Cornell Aeronautical Laboratory by Frank Rosenblatt
- Implemented in custom-built hardware as the “Mark 1 perceptron,” designed for image recognition
- 400 photocells, randomly connected to the “neurons.” Weights were encoded in potentiometers
- Weight updates during learning were performed by electric motors.

Perceptron in the News...

NEW NAVY DEVICE LEARNS BY DOING

Psychologist Shows Embryo of Computer Designed to Read and Grow Wiser

WASHINGTON, July 7 (UPI)—The Navy revealed the embryo of an electronic computer today that it expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence.

The embryo—the Weather Bureau's \$2,000,000 "704" computer—learned to differentiate between right and left after fifty attempts in the Navy's demonstration for newsmen.

The service said it would use this principle to build the first of its Perceptron thinking machines that will be able to read and write. It is expected to be finished in about a year at a cost of \$100,000.

Dr. Frank Rosenblatt, designer of the Perceptron, conducted the demonstration. He said the machine would be the first device to think as the human brain. As do human be-

ings, Perceptron will make mistakes at first, but will grow wiser as it gains experience, he said.

Dr. Rosenblatt, a research psychologist at the Cornell Aeronautical Laboratory, Buffalo, said Perceptrons might be fired to the planets as mechanical space explorers.

Without Human Controls

The Navy said the perceptron would be the first non-living mechanism "capable of receiving, recognizing and identifying its surroundings without any human training or control."

The "brain" is designed to remember images and information it has perceived itself. Ordinary computers remember only what is fed into them on punch cards or magnetic tape.

Later Perceptrons will be able to recognize people and call out their names and instantly translate speech in one language to speech or writing in another language, it was predicted.

Mr. Rosenblatt said in principle it would be possible to build brains that could reproduce themselves on an assembly line and which would be conscious of their existence.

1958 New York Times...

In today's demonstration, the "704" was fed two cards, one with squares marked on the left side and the other with squares on the right side.

Learns by Doing

In the first fifty trials, the machine made no distinction between them. It then started registering a "Q" for the left squares and "O" for the right squares.

Dr. Rosenblatt said he could explain why the machine learned only in highly technical terms. But he said the computer had undergone a "self-induced change in the wiring diagram."

The first Perceptron will have about 1,000 electronic "association cells" receiving electrical impulses from an eye-like scanning device with 400 photo-cells. The human brain has 10,000,000,000 responsive cells, including 100,000,000 connections with the eyes.

Perceptron in the News...

NEW NAVY DEVICE LEARNS BY DOING

Psychologist Shows Embryo of Computer Designed to Read and Grow Wiser

WASHINGTON, July 7 (UPI)—The Navy revealed the embryo of an electronic computer today that it expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence.

The embryo—the Weather Bureau's \$2,000,000 "704" computer—learned to differentiate between right and left after fifty attempts in the Navy's demonstration for newsmen.

The service said it would use this principle to build the first of its Perceptron thinking machines that will be able to read and write. It is expected to be finished in about a year at a cost of \$100,000.

Mr. Frank Rosenblatt, designer of the Perceptron, conducted the demonstration. He said the machine would be the first device to think as the human brain. As do human be-

ings, Perceptron will make mistakes at first, but will grow wiser as it gains experience, he said.

Dr. Rosenblatt, a research psychologist at the Cornell Aeronautical Laboratory, Buffalo, said Perceptrons might be fired to the planets as mechanical space explorers.

Without Human Controls

The Navy said the perceptron would be the first non-living mechanism "capable of receiving, recognizing and identifying its surroundings without any human training or control."

The "brain" is designed to remember images and information it has perceived itself. Ordinary computers remember only what is fed into them on punch cards or magnetic tape.

Later Perceptrons will be able to recognize people and call out their names and instantly translate speech in one language to speech or writing in another language, it was predicted.

Mr. Rosenblatt said in principle it would be possible to build brains that could reproduce themselves on an assembly line and which would be conscious of their existence.

1958 New York Times...

In today's demonstration, the "704" was fed two cards, one with squares marked on the left side and the other with squares on the right side.

Learns by Doing

In the first fifty trials, the machine made no distinction between them. It then started registering a "Q" for the left squares and "O" for the right squares.

Dr. Rosenblatt said he could explain why the machine learned only in highly technical terms. But he said the computer had undergone a "self-induced change in the wiring diagram."

The first Perceptron will have about 1,000 electronic "association cells" receiving electrical impulses from an eye-like scanning device with 400 photo-cells. The human brain has 10,000,000,000 responsive cells, including 100,000,000 connections with the eyes.

Perceptron Algorithm

- **Online** algorithm: process one data point at each round
 - Take x_i ; apply the current model to make a prediction for it
 - If prediction is **correct**, proceed
 - **Else**, correct model: add feature vector w.r.t. correct output & subtract feature vector w.r.t. predicted (wrong) output

Perceptron Algorithm

input: labeled data \mathcal{D}

initialize $w^{(0)} = \mathbf{0}$

initialize $k = 0$ (**number of mistakes**)

repeat

 get new training example (x_i, y_i)

 predict $\hat{y}_i = \arg \max_{y \in \mathcal{Y}} w^{(k)} \cdot \phi(x_i, y)$

if $\hat{y}_i \neq y_i$ **then**

 update $w^{(k+1)} = w^{(k)} + \phi(x_i, y_i) - \phi(x_i, \hat{y}_i)$

 increment k

end if

until maximum number of epochs

output: model weights w

Perceptron's Mistake Bound

A couple definitions:

- the training data is **linearly separable** with margin $\gamma > 0$ iff there is a weight vector u with $\|u\| = 1$ such that

$$u \cdot \phi(x_i, y_i) \geq u \cdot \phi(x_i, y'_i) + \gamma, \quad \forall i, \quad \forall y'_i \neq y_i.$$

- **radius** of the data: $R = \max_{i, y'_i \neq y_i} \|\phi(x_i, y_i) - \phi(x_i, y'_i)\|$.

Perceptron's Mistake Bound

A couple definitions:

- the training data is **linearly separable** with margin $\gamma > 0$ iff there is a weight vector u with $\|u\| = 1$ such that

$$u \cdot \phi(x_i, y_i) \geq u \cdot \phi(x_i, y'_i) + \gamma, \quad \forall i, \quad \forall y'_i \neq y_i.$$

- **radius** of the data: $R = \max_{i, y'_i \neq y_i} \|\phi(x_i, y_i) - \phi(x_i, y'_i)\|$.

Then we have the following bound of the **number of mistakes**:

Theorem (Novikoff (1962))

The perceptron algorithm is guaranteed to find a separating hyperplane after at most $\frac{R^2}{\gamma^2}$ mistakes.

One-Slide Proof

- Lower bound on $\|w^{(k+1)}\|$:

$$\begin{aligned} \mathbf{u} \cdot \mathbf{w}^{(k+1)} &= \mathbf{u} \cdot \mathbf{w}^{(k)} + \mathbf{u} \cdot (\phi(\mathbf{x}_i, \mathbf{y}_i) - \phi(\mathbf{x}_i, \hat{\mathbf{y}}_i)) \\ &\geq \mathbf{u} \cdot \mathbf{w}^{(k)} + \gamma \\ &\geq k\gamma. \end{aligned}$$

Hence $\|\mathbf{w}^{(k+1)}\| = \|\mathbf{u}\| \cdot \|\mathbf{w}^{(k+1)}\| \geq \mathbf{u} \cdot \mathbf{w}^{(k+1)} \geq k\gamma$ (from CSI).

One-Slide Proof

- Lower bound on $\|w^{(k+1)}\|$:

$$\begin{aligned} \mathbf{u} \cdot \mathbf{w}^{(k+1)} &= \mathbf{u} \cdot \mathbf{w}^{(k)} + \mathbf{u} \cdot (\phi(\mathbf{x}_i, \mathbf{y}_i) - \phi(\mathbf{x}_i, \hat{\mathbf{y}}_i)) \\ &\geq \mathbf{u} \cdot \mathbf{w}^{(k)} + \gamma \\ &\geq k\gamma. \end{aligned}$$

Hence $\|\mathbf{w}^{(k+1)}\| = \|\mathbf{u}\| \cdot \|\mathbf{w}^{(k+1)}\| \geq \mathbf{u} \cdot \mathbf{w}^{(k+1)} \geq k\gamma$ (from CSI).

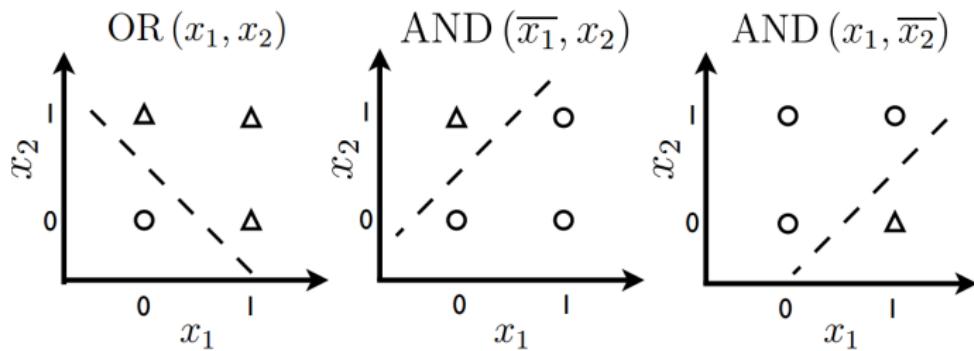
- Upper bound on $\|w^{(k+1)}\|$:

$$\begin{aligned} \|\mathbf{w}^{(k+1)}\|^2 &= \|\mathbf{w}^{(k)}\|^2 + \|\phi(\mathbf{x}_i, \mathbf{y}_i) - \phi(\mathbf{x}_i, \hat{\mathbf{y}}_i)\|^2 \\ &\quad + 2\mathbf{w}^{(k)} \cdot (\phi(\mathbf{x}_i, \mathbf{y}_i) - \phi(\mathbf{x}_i, \hat{\mathbf{y}}_i)) \\ &\leq \|\mathbf{w}^{(k)}\|^2 + R^2 \\ &\leq kR^2. \end{aligned}$$

Equating both sides, we get $(k\gamma)^2 \leq kR^2 \Rightarrow k \leq R^2/\gamma^2$ (QED).

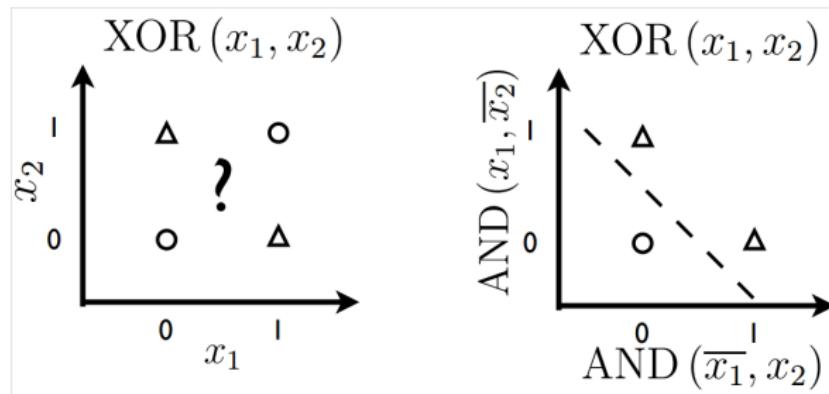
What a Simple Perceptron Can and Can't Do

- Remember: the decision boundary is linear (**linear classifier**)
- It **can** solve linearly separable problems (OR, AND)



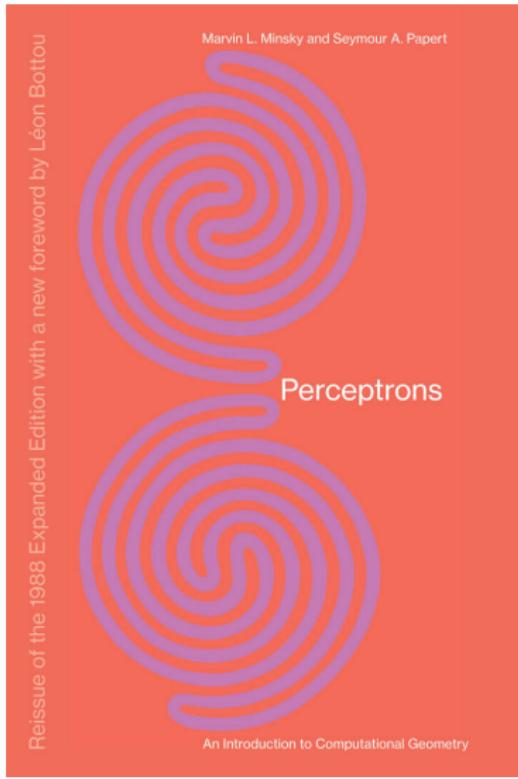
What a Simple Perceptron Can and Can't Do

- ... but it **can't** solve **non-linearly separable** problems such as simple XOR (unless input is transformed into a better representation):



- This result is often attributed to Minsky and Papert (1969) but was known well before.

Limitations of the Perceptron



Minsky and Papert (1969):

- Shows limitations of multi-layer perceptrons and fostered an “AI winter” period.

More tomorrow at Bhiksha’s lecture!

Outline

① Data and Feature Representation

② Perceptron

③ Naive Bayes

④ Logistic Regression

⑤ Support Vector Machines

⑥ Regularization

⑦ Non-Linear Classifiers

Probabilistic Models

- For a moment, forget linear classifiers and parameter vectors w
- Let's assume our goal is to model the conditional probability of output labels y given inputs x , i.e. $P(y|x)$
- If we can define this distribution, then classification becomes:

$$\hat{y} = \arg \max_{y \in \mathcal{Y}} P(y|x)$$

Bayes Rule

- One way to model $P(y|x)$ is through Bayes Rule:

$$P(y|x) = \frac{P(y)P(x|y)}{P(x)}$$

$$\arg \max_y P(y|x) = \arg \max_y P(y)P(x|y)$$

(since x is fixed!)

- $P(y)P(x|y) = P(x,y)$: a joint probability
- Modeling the joint input-output distribution is at the core of generative models
 - Because we model a distribution that can randomly generate outputs *and* inputs, not just outputs

Naive Bayes

Assume that an input x is partitioned as v_1, \dots, v_L , where $v_k \in \mathcal{V}_k$

Example:

- x is a document of length L
- v_k is the k^{th} token (a word)
- The set $\mathcal{V}_k = \mathcal{V}$ is a fixed vocabulary (all tokens drawn from \mathcal{V})

Naive Bayes Assumption
(conditional independence)

$$P(\underbrace{v_1, \dots, v_L}_{x} | y) = \prod_{k=1}^L P(v_k | y)$$

Multinomial Naive Bayes

$$P(x, y) = P(y) P(\underbrace{v_1, \dots, v_L}_x | y) = P(y) \prod_{k=1}^L P(v_k | y)$$

- All tokens are conditionally independently, given the topic
- The word order doesn't change $P(x, y)$ (bag-of-words assumption)

Small caveat: we assumed that the document has a fixed length L .

This is not realistic.

How to deal with variable length?

Multinomial Naive Bayes – Arbitrary Length

Solution: introduce a distribution over document length $P(|\mathbf{x}|)$

- e.g. a Poisson distribution.

We get:

$$P(\mathbf{x}, \mathbf{y}) = P(\mathbf{y}) P(|\mathbf{x}|) \underbrace{\prod_{k=1}^{|x|} P(v_k | \mathbf{y})}_{P(\mathbf{x}|\mathbf{y})}$$

$P(|\mathbf{x}|)$ is constant (independent of \mathbf{y}), so nothing really changes

- the posterior $P(\mathbf{y}|\mathbf{x})$ is the same as before.

What Does This Buy Us?

$$P(\underbrace{v_1, \dots, v_L}_x | \mathbf{y}) = \prod_{k=1}^L P(v_k | \mathbf{y})$$

What do we gain with the Naive Bayes assumption?

What Does This Buy Us?

$$P(\underbrace{v_1, \dots, v_L}_x | \mathbf{y}) = \prod_{k=1}^L P(v_k | \mathbf{y})$$

What do we gain with the Naive Bayes assumption?

- A huge reduction in the number of parameters!
- If we haven't done any factorization assumption, how many parameters would be required for expressing $P(v_1, \dots, v_L | \mathbf{y})$?

What Does This Buy Us?

$$P(\underbrace{v_1, \dots, v_L}_x | \mathbf{y}) = \prod_{k=1}^L P(v_k | \mathbf{y})$$

What do we gain with the Naive Bayes assumption?

- A huge reduction in the number of parameters!
- If we haven't done any factorization assumption, how many parameters would be required for expressing $P(v_1, \dots, v_L | \mathbf{y})$? $O(|\mathcal{V}|^L)$
- And how many parameters with Naive Bayes?

What Does This Buy Us?

$$P(\underbrace{v_1, \dots, v_L}_x | \mathbf{y}) = \prod_{k=1}^L P(v_k | \mathbf{y})$$

What do we gain with the Naive Bayes assumption?

- A huge reduction in the number of parameters!
- If we haven't done any factorization assumption, how many parameters would be required for expressing $P(v_1, \dots, v_L | \mathbf{y})$? $O(|\mathcal{V}|^L)$
- And how many parameters with Naive Bayes? $O(|\mathcal{V}|)$

What Does This Buy Us?

$$P(\underbrace{v_1, \dots, v_L}_x | \mathbf{y}) = \prod_{k=1}^L P(v_k | \mathbf{y})$$

What do we gain with the Naive Bayes assumption?

- A huge reduction in the number of parameters!
- If we haven't done any factorization assumption, how many parameters would be required for expressing $P(v_1, \dots, v_L | \mathbf{y})$? $O(|\mathcal{V}|^L)$
- And how many parameters with Naive Bayes? $O(|\mathcal{V}|)$

Less parameters \implies Less computation; less risk of overfitting

(Though we may underfit if our independence assumptions are too strong.)

Naive Bayes – Learning

$$P(\mathbf{y})P(\underbrace{v_1, \dots, v_L}_{\mathbf{x}} | \mathbf{y}) = P(\mathbf{y}) \prod_{k=1}^L P(v_k | \mathbf{y})$$

- Input: dataset $\mathcal{D} = \{(\mathbf{x}_t, \mathbf{y}_t)\}_{t=1}^N$ (**examples assumed i.i.d.**)
- Parameters $\Theta = \{P(\mathbf{y}), P(v|\mathbf{y})\}$
- **Objective: Maximum Likelihood Estimation (MLE):** choose parameters that maximize the likelihood of observed data

$$\mathcal{L}(\Theta; \mathcal{D}) = \prod_{t=1}^N P(\mathbf{x}_t, \mathbf{y}_t) = \prod_{t=1}^N \left(P(\mathbf{y}_t) \prod_{k=1}^L P(v_k(\mathbf{x}_t) | \mathbf{y}_t) \right)$$

$$\hat{\Theta} = \arg \max_{\Theta} \prod_{t=1}^N \left(P(\mathbf{y}_t) \prod_{k=1}^L P(v_k(\mathbf{x}_t) | \mathbf{y}_t) \right)$$

Naive Bayes – Learning via MLE

For the multinomial Naive Bayes model, MLE has a **closed form solution!!**
It all boils down to counting and normalizing!!
(The proof is left as an exercise...)

Naive Bayes – Learning via MLE

$$\hat{\Theta} = \arg \max_{\Theta} \prod_{t=1}^N \left(P(y_t) \prod_{k=1}^L P(v_k(x_t) | y_t) \right)$$

$$\hat{P}(y) = \frac{\sum_{t=1}^N [[y_t = y]]}{N}$$

$$\hat{P}(v|y) = \frac{\sum_{t=1}^N \sum_{k=1}^L [[v_k(x_t) = v \text{ and } y_t = y]]}{L \sum_{t=1}^N [[y_t = y]]}$$

$[[X]]$ is 1 if property X holds, 0 otherwise (Iverson notation)
Fraction of times a feature appears in training cases of a given label

Naive Bayes Example

- Corpus of movie reviews: 7 examples for **training**

Doc	Words	Class
1	Great movie, excellent plot, renown actors	Positive
2	I had not seen a fantastic plot like this in good 5 years. Amazing!!!	Positive
3	Lovely plot, amazing cast, somehow I am in love with the bad guy	Positive
4	Bad movie with great cast, but very poor plot and unimaginative ending	Negative
5	I hate this film, it has nothing original	Negative
6	Great movie, but not...	Negative
7	Very bad movie, I have no words to express how I dislike it	Negative

Naive Bayes Example

- **Features:** adjectives (bag-of-words)

Doc	Words	Class
1	Great movie, excellent plot, renowned actors	Positive
2	I had not seen a fantastic plot like this in good 5 years. amazing !!!	Positive
3	Lovely plot, amazing cast, somehow I am in love with the bad guy	Positive
4	Bad movie with great cast, but very poor plot and unimaginative ending	Negative
5	I hate this film, it has nothing original. Really bad	Negative
6	Great movie, but not...	Negative
7	Very bad movie, I have no words to express how I dislike it	Negative

Naive Bayes Example

Relative frequency:

Priors:

$$P(\text{positive}) = \frac{\sum_{t=1}^N [[y_t = \text{positive}]]}{N} = 3/7 = 0.43$$

$$P(\text{negative}) = \frac{\sum_{t=1}^N [[y_t = \text{negative}]]}{N} = 4/7 = 0.57$$

Assume standard pre-processing: tokenization, lowercasing, punctuation removal (except special punctuation like !!!)

Naive Bayes Example

Likelihoods: Count adjective v in class y / adjectives in y

$$\hat{P}(v|y) = \frac{\sum_{t=1}^N \sum_{k=1}^L [[v_k(x_t) = v \text{ and } y_t = y]]}{L \sum_{t=1}^N [[y_t = y]]}$$

$P(\text{amazing} \text{positive})$	= 2/10	$P(\text{amazing} \text{negative})$	= 0/8
$P(\text{bad} \text{positive})$	= 1/10	$P(\text{bad} \text{negative})$	= 3/8
$P(\text{excellent} \text{positive})$	= 1/10	$P(\text{excellent} \text{negative})$	= 0/8
$P(\text{fantastic} \text{positive})$	= 1/10	$P(\text{fantastic} \text{negative})$	= 0/8
$P(\text{good} \text{positive})$	= 1/10	$P(\text{good} \text{negative})$	= 0/8
$P(\text{great} \text{positive})$	= 1/10	$P(\text{great} \text{negative})$	= 2/8
$P(\text{lovely} \text{positive})$	= 1/10	$P(\text{lovely} \text{negative})$	= 0/8
$P(\text{original} \text{positive})$	= 0/10	$P(\text{original} \text{negative})$	= 1/8
$P(\text{poor} \text{positive})$	= 0/10	$P(\text{poor} \text{negative})$	= 1/8
$P(\text{renowned} \text{positive})$	= 1/10	$P(\text{renowned} \text{negative})$	= 0/8
$P(\text{unimaginative} \text{positive})$	= 0/10	$P(\text{unimaginative} \text{negative})$	= 1/8

Naive Bayes Example

Given a new segment to classify (**test time**):

Doc	Words	Class
8	This was a fantastic story, good , lovely	???

Final decision

$$\hat{y} = \arg \max_y \left(P(y) \prod_{k=1}^L P(v_k|y) \right)$$

$$P(\text{positive}) * P(\text{fantastic}|\text{positive}) * P(\text{good}|\text{positive}) * P(\text{lovely}|\text{positive})$$

$$3/7 * 1/10 * 1/10 * 1/10 = 0.00043$$

$$P(\text{negative}) * P(\text{fantastic}|\text{negative}) * P(\text{good}|\text{negative}) * P(\text{lovely}|\text{negative})$$

$$4/7 * 0/8 * 0/8 * 0/8 = 0$$

So: *sentiment = positive*

Naive Bayes Example

Given a new segment to classify (**test time**):

Doc	Words	Class
9	Great plot, great cast, great everything	???

Final decision

$$P(\text{positive}) * P(\text{great}|\text{positive}) * P(\text{great}|\text{positive}) * P(\text{great}|\text{positive})$$

$$3/7 * 1/10 * 1/10 * 1/10 = 0.00043$$

$$P(\text{negative}) * P(\text{great}|\text{negative}) * P(\text{great}|\text{negative}) * P(\text{great}|\text{negative})$$

$$4/7 * 2/8 * 2/8 * 2/8 = 0.00893$$

So: *sentiment = negative*

Naive Bayes Example

But if the new segment to classify (**test time**) is:

Doc	Words	Class
10	Boring movie, annoying plot, unimaginative ending	???

Final decision

$$P(\text{positive}) * P(\text{boring}|\text{positive}) * P(\text{annoying}|\text{positive}) * P(\text{unimaginative}|\text{positive})$$

$$3/7 * 0/10 * 0/10 * 0/10 = 0$$

$$P(\text{negative}) * P(\text{boring}|\text{negative}) * P(\text{annoying}|\text{negative}) * P(\text{unimaginative}|\text{negative})$$

$$4/7 * 0/8 * 0/8 * 1/8 = 0$$

So: *sentiment* = ???

Laplace Smoothing

Add smoothing to feature counts (add 1 to every count):

$$\hat{P}(v|y) = \frac{\sum_{t=1}^N \sum_{k=1}^L [[v_k(x_t) = v \text{ and } y_t = y]] + 1}{L \sum_{t=1}^N [[y_t = y]] + |\mathcal{V}|}$$

where $|\mathcal{V}|$ = number of distinct adjectives in training (all classes) = 12

Doc	Words	Class
11	Boring movie, annoying plot, unimaginative ending	???

Final decision

$$P(\text{positive}) * P(\text{boring}|\text{positive}) * P(\text{annoying}|\text{positive}) * P(\text{unimaginative}|\text{positive})$$

$$3/7 * ((0 + 1)/(10 + 12)) * ((0 + 1)/(10 + 12)) * ((0 + 1)/(10 + 12)) = 0.000040$$

$$P(\text{negative}) * P(\text{boring}|\text{negative}) * P(\text{annoying}|\text{negative}) * P(\text{unimaginative}|\text{negative})$$

$$4/7 * ((0 + 1)/(8 + 12)) * ((0 + 1)/(8 + 12)) * ((1 + 1)/(8 + 12)) = 0.000143$$

So: *sentiment = negative*



Finally...

Multinomial Naive Bayes is a Linear Classifier!

One Slide Proof

- Let $b_y = \log P(y)$, $\forall y \in \mathcal{Y}$
- Let $[w_y]_v = \log P(v|y)$, $\forall y \in \mathcal{Y}, v \in \mathcal{V}$
- Let $[\psi(x)]_v = \sum_{k=1}^L [[v_k(x) = v]]$, $\forall v \in \mathcal{V}$ ($\#$ times v occurs in x)

$$\begin{aligned}\arg \max_y P(y|x) &\propto \arg \max_y \left(P(y) \prod_{k=1}^L P(v_k(x)|y) \right) \\&= \arg \max_y \left(\log P(y) + \sum_{k=1}^L \log P(v_k(x)|y) \right) \\&= \arg \max_y \left(\underbrace{\log P(y)}_{b_y} + \sum_{v \in \mathcal{V}} [\psi(x)]_v \underbrace{\log P(v|y)}_{[w_y]_v} \right) \\&= \arg \max_y (w_y \cdot \psi(x) + b_y).\end{aligned}$$

Discriminative versus Generative

- Generative models attempt to model inputs and outputs
 - e.g., Naive Bayes = MLE of joint distribution $P(x, y)$
 - Statistical model must explain generation of input
 - Can we sample a document from the multinomial Naive Bayes model?
How?

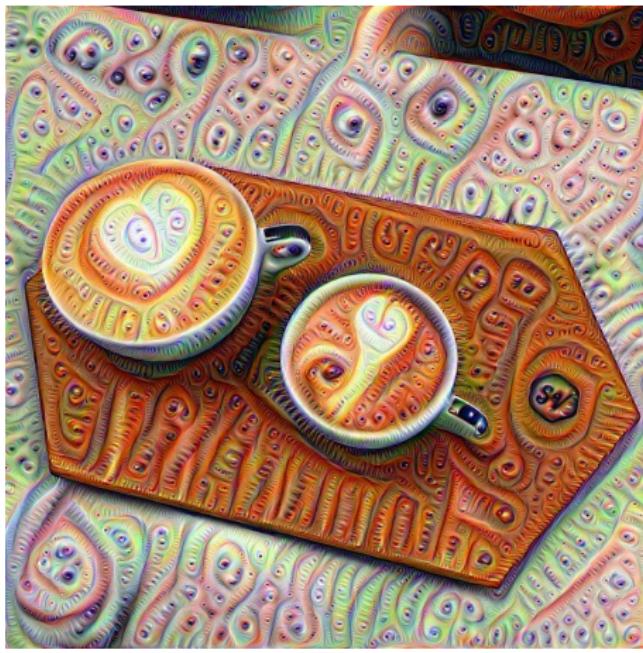
Discriminative versus Generative

- Generative models attempt to model inputs and outputs
 - e.g., Naive Bayes = MLE of joint distribution $P(\mathbf{x}, \mathbf{y})$
 - Statistical model must explain generation of input
 - Can we sample a document from the multinomial Naive Bayes model?
How?
- Occam's Razor: why model input?
- Discriminative models
 - Use loss function that directly optimizes $P(\mathbf{y}|\mathbf{x})$ (or something related)
 - Logistic Regression – MLE of $P(\mathbf{y}|\mathbf{x})$
 - Perceptron and SVMs – minimize classification error

Discriminative versus Generative

- Generative models attempt to model inputs and outputs
 - e.g., Naive Bayes = MLE of joint distribution $P(\mathbf{x}, \mathbf{y})$
 - Statistical model must explain generation of input
 - Can we sample a document from the multinomial Naive Bayes model? How?
- Occam's Razor: why model input?
- Discriminative models
 - Use loss function that directly optimizes $P(\mathbf{y}|\mathbf{x})$ (or something related)
 - Logistic Regression – MLE of $P(\mathbf{y}|\mathbf{x})$
 - Perceptron and SVMs – minimize classification error
- Generative and discriminative models use $P(\mathbf{y}|\mathbf{x})$ for prediction
 - They differ only on what distribution they use to set \mathbf{w}

Coffee-break!



So far

We have covered:

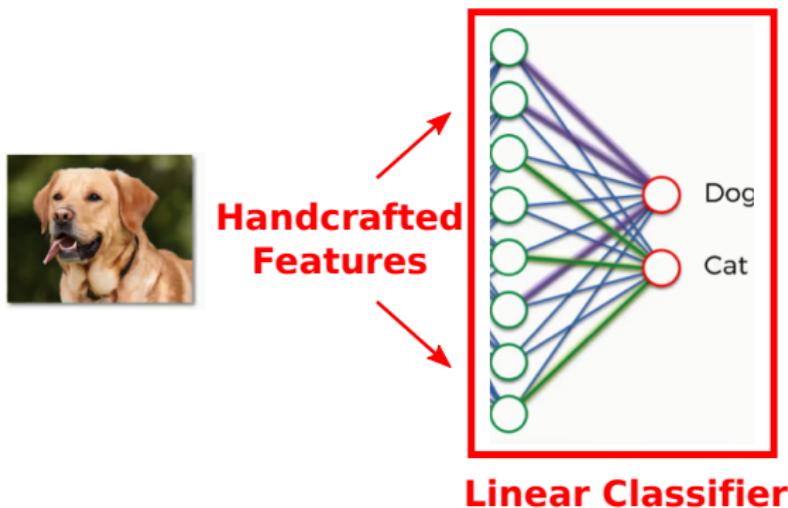
- The perceptron algorithm
- (Multinomial) Naive Bayes.

We saw that both are instances of **linear classifiers**.

Perceptron finds a separating hyperplane (if it exists), Naive Bayes is a generative probabilistic model

Next: a **discriminative** probabilistic model.

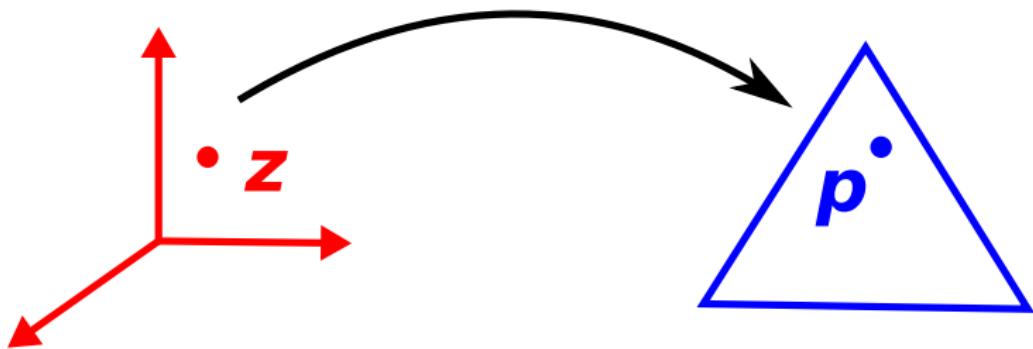
Reminder



$$\hat{y} = \text{argmax} (\mathbf{W}\psi(x) + \mathbf{b}), \quad \mathbf{W} = \begin{bmatrix} \vdots \\ w_y^\top \\ \vdots \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} \vdots \\ b_y \\ \vdots \end{bmatrix}.$$

Key Problem

How to map from a set of label scores $\mathbb{R}^{|\mathcal{Y}|}$ to a probability distribution over \mathcal{Y} ?



We'll see two mappings: softmax (next) and sparsemax (later).

Outline

① Data and Feature Representation

② Perceptron

③ Naive Bayes

④ Logistic Regression

⑤ Support Vector Machines

⑥ Regularization

⑦ Non-Linear Classifiers

Logistic Regression

Define a conditional probability:

$$P(y|x) = \frac{\exp(w \cdot \phi(x, y))}{Z_x}, \quad \text{where } Z_x = \sum_{y' \in \mathcal{Y}} \exp(w \cdot \phi(x, y'))$$

This operation (exponentiating and normalizing) is called the **softmax transformation** (more later!)

Note: still a linear classifier

$$\begin{aligned} \arg \max_y P(y|x) &= \arg \max_y \frac{\exp(w \cdot \phi(x, y))}{Z_x} \\ &= \arg \max_y \exp(w \cdot \phi(x, y)) \\ &= \arg \max_y w \cdot \phi(x, y) \end{aligned}$$

Binary Logistic Regression

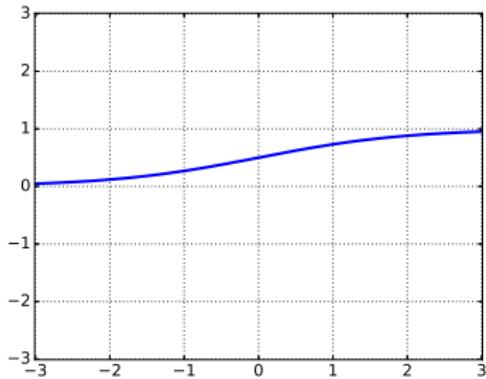
Binary labels ($\mathcal{Y} = \{\pm 1\}$)

$$\begin{aligned} P(y = +1|x) &= \frac{\exp(\mathbf{v} \cdot \psi(x) + c)}{1 + \exp(\mathbf{v} \cdot \psi(x) + c)} \\ &= \frac{1}{1 + \exp(-\mathbf{v} \cdot \psi(x) - c)} \\ &= \sigma(\mathbf{v} \cdot \psi(x) + c). \end{aligned}$$

This is called a **sigmoid transformation** (more later!)

Sigmoid Transformation

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



- Widely used in neural networks (wait for tomorrow!)
- Can be regarded as a 2D softmax
- “Squashes” a real number between 0 and 1
- The output can be interpreted as a probability
- Positive, bounded, strictly increasing

Multinomial Logistic Regression

$$P_w(y|x) = \frac{\exp(w \cdot \phi(x, y))}{Z_x}$$

- How do we learn weights w ?
- Set w to maximize the **conditional log-likelihood** of training data:

$$\begin{aligned}\hat{w} &= \arg \max_{w \in \mathbb{R}^D} \log \left(\prod_{t=1}^N P_w(y_t|x_t) \right) = \arg \min_{w \in \mathbb{R}^D} - \sum_{t=1}^N \log P_w(y_t|x_t) = \\ &= \arg \min_{w \in \mathbb{R}^D} \sum_{t=1}^N \left(\log \sum_{y'_t} \exp(w \cdot \phi(x_t, y'_t)) - w \cdot \phi(x_t, y_t) \right),\end{aligned}$$

i.e., set w to assign as much probability mass as possible to the correct labels!

Logistic Regression

- This objective function is **convex**
- Therefore any local minimum is a global minimum
- No closed form solution, but lots of numerical techniques
 - Gradient methods (gradient descent, conjugate gradient)
 - Quasi-Newton methods (L-BFGS, ...)

Logistic Regression

- This objective function is **convex**
- Therefore any local minimum is a global minimum
- No closed form solution, but lots of numerical techniques
 - Gradient methods (gradient descent, conjugate gradient)
 - Quasi-Newton methods (L-BFGS, ...)
- Logistic Regression = **Maximum Entropy**: maximize entropy subject to constraints on features
- Proof left as an exercise!

Recap: Convex functions

Pro: Guarantee of a global minima ✓



Figure: Illustration of a convex function. The line segment between any two points on the graph lies entirely above the curve.

Recap: Iterative Descent Methods

Goal: find the minimum/minimizer of $f : \mathbb{R}^d \rightarrow \mathbb{R}$

- Proceed in **small steps** in the **optimal direction** till a **stopping criterion** is met.
- **Gradient descent:** updates of the form: $x^{(k+1)} \leftarrow x^{(k)} - \eta_k \nabla f(x^{(k)})$

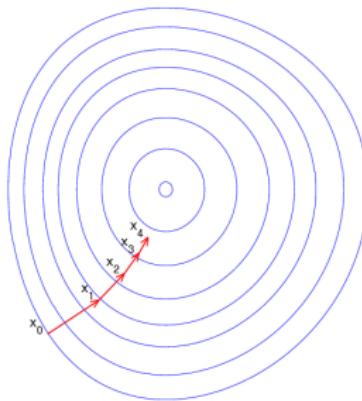


Figure: Illustration of gradient descent. The red lines correspond to steps taken in the negative gradient direction.

Gradient Descent

- Let $L(\mathbf{w}; (x, y)) = \log \sum_{y'} \exp(\mathbf{w} \cdot \phi(x, y')) - \mathbf{w} \cdot \phi(x, y)$
- This is our **loss function!**
- We want to find $\arg \min_{\mathbf{w}} \sum_{t=1}^N L(\mathbf{w}; (x_t, y_t))$
 - Set $\mathbf{w}^0 = \mathbf{0}$
 - Iterate until convergence (for suitable stepsize η_k):

$$\begin{aligned}\mathbf{w}^{k+1} &= \mathbf{w}^k - \eta_k \nabla_{\mathbf{w}} \left(\sum_{t=1}^N L(\mathbf{w}; (x_t, y_t)) \right) \\ &= \mathbf{w}^k - \eta_k \sum_{t=1}^N \nabla_{\mathbf{w}} L(\mathbf{w}; (x_t, y_t))\end{aligned}$$

- $\nabla_{\mathbf{w}} L(\mathbf{w})$ is gradient of L w.r.t. \mathbf{w}
- Gradient descent will always find the optimal \mathbf{w}

Stochastic Gradient Descent

If the dataset is large, we'd better do SGD instead, for more frequent updates:

- Set $\mathbf{w}^0 = \mathbf{0}$
- Iterate until convergence
 - Pick (\mathbf{x}_t, y_t) randomly
 - Update $\mathbf{w}^{k+1} = \mathbf{w}^k - \eta_k \nabla_{\mathbf{w}} L(\mathbf{w}; (\mathbf{x}_t, y_t))$
- i.e. we approximate the true gradient with a noisy, unbiased, gradient, based on **a single sample**
- Variants exist in-between (mini-batches)
- All guaranteed to find the optimal \mathbf{w} (for suitable step sizes)

Computing the Gradient

- For this to work, we need to be able to compute $\nabla_{\mathbf{w}} L(\mathbf{w}; (\mathbf{x}_t, \mathbf{y}_t))$, where

$$L(\mathbf{w}; (\mathbf{x}, \mathbf{y})) = \log \sum_{\mathbf{y}'} \exp(\mathbf{w} \cdot \phi(\mathbf{x}, \mathbf{y}')) - \mathbf{w} \cdot \phi(\mathbf{x}, \mathbf{y})$$

Some reminders:

- $\nabla_{\mathbf{w}} \log F(\mathbf{w}) = \frac{1}{F(\mathbf{w})} \nabla_{\mathbf{w}} F(\mathbf{w})$
- $\nabla_{\mathbf{w}} \exp F(\mathbf{w}) = \exp(F(\mathbf{w})) \nabla_{\mathbf{w}} F(\mathbf{w})$

Computing the Gradient

$$\begin{aligned}\nabla_w L(w; (x, y)) &= \nabla_w \left(\log \sum_{y'} \exp(w \cdot \phi(x, y')) - w \cdot \phi(x, y) \right) \\&= \nabla_w \log \sum_{y'} \exp(w \cdot \phi(x, y')) - \nabla_w w \cdot \phi(x, y) \\&= \frac{1}{\sum_{y'} \exp(w \cdot \phi(x, y'))} \sum_{y'} \nabla_w \exp(w \cdot \phi(x, y')) - \phi(x, y) \\&= \frac{1}{Z_x} \sum_{y'} \exp(w \cdot \phi(x, y')) \nabla_w w \cdot \phi(x, y') - \phi(x, y) \\&= \sum_{y'} \frac{\exp(w \cdot \phi(x, y'))}{Z_x} \phi(x, y') - \phi(x, y) \\&= \sum_{y'} P_w(y' | x) \phi(x, y') - \phi(x, y).\end{aligned}$$

The gradient equals the “difference between the **expected features under the current model** and the **true features**.”

Logistic Regression Summary

- Define conditional probability

$$P_w(y|x) = \frac{\exp(w \cdot \phi(x, y))}{Z_x}$$

- Set weights to maximize conditional log-likelihood of training data:

$$w = \arg \max_w \sum_t \log P_w(y_t|x_t) = \arg \min_w \sum_t L(w; (x_t, y_t))$$

- Can find the gradient and run gradient descent (or any gradient-based optimization algorithm)

$$\nabla_w L(w; (x, y)) = \sum_{y'} P_w(y'|x) \phi(x, y') - \phi(x, y)$$

The Story So Far

- Naive Bayes is **generative**: maximizes **joint** likelihood
 - closed form solution (boils down to **counting and normalizing**)
- Logistic regression is **discriminative**: maximizes **conditional** likelihood
 - also called log-linear model and max-entropy classifier
 - no closed form solution
 - stochastic gradient updates look like

$$\mathbf{w}^{k+1} = \mathbf{w}^k + \eta \left(\phi(\mathbf{x}, \mathbf{y}) - \sum_{\mathbf{y}'} P_{\mathbf{w}}(\mathbf{y}'|\mathbf{x})\phi(\mathbf{x}, \mathbf{y}') \right)$$

- Perceptron is a discriminative, non-probabilistic classifier
 - perceptron's updates look like

$$\mathbf{w}^{k+1} = \mathbf{w}^k + \phi(\mathbf{x}, \mathbf{y}) - \phi(\mathbf{x}, \hat{\mathbf{y}})$$

SGD updates for logistic regression and perceptron's updates look similar!

Maximizing Margin

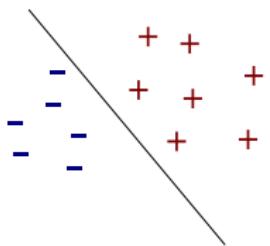
- For a training set \mathcal{D}
- Margin of a weight vector w is smallest γ such that

$$w \cdot \phi(x_t, y_t) - w \cdot \phi(x_t, y') \geq \gamma$$

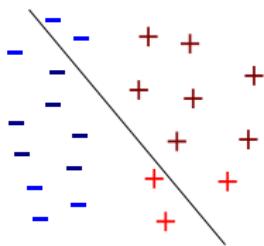
- for every training instance $(x_t, y_t) \in \mathcal{D}, y' \in \mathcal{Y}$

Margin

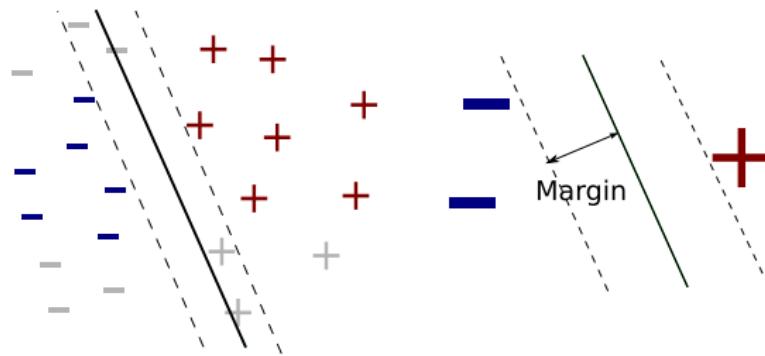
Training



Testing



Denote the value of the margin by γ



Maximizing Margin

- Intuitively maximizing margin makes sense
- More importantly, generalization error to unseen test data is proportional to the inverse of the margin

$$\epsilon \propto \frac{R^2}{\gamma^2 \times N}$$

- **Perceptron:**

- If a training set is separable by some margin, the perceptron will find a w that separates the data
- However, the perceptron does not pick w to maximize the margin!

Outline

① Data and Feature Representation

② Perceptron

③ Naive Bayes

④ Logistic Regression

⑤ Support Vector Machines

⑥ Regularization

⑦ Non-Linear Classifiers

Maximizing Margin

Let $\gamma > 0$

$$\max_{\|\mathbf{w}\| \leq 1} \gamma$$

such that:

$$\mathbf{w} \cdot \phi(\mathbf{x}_t, \mathbf{y}_t) - \mathbf{w} \cdot \phi(\mathbf{x}_t, \mathbf{y}') \geq \gamma$$

$$\forall (\mathbf{x}_t, \mathbf{y}_t) \in \mathcal{D}$$

$$\text{and } \mathbf{y}' \in \mathcal{Y}$$

- Note: algorithm still **minimizes error** if data is separable
- $\|\mathbf{w}\|$ is bound since scaling trivially produces larger margin

Max Margin = Min Norm

Let $\gamma > 0$

Max Margin:

$$\max_{\|\mathbf{w}\| \leq 1} \gamma$$

such that:

=

$$\mathbf{w} \cdot \phi(\mathbf{x}_t, \mathbf{y}_t) - \mathbf{w} \cdot \phi(\mathbf{x}_t, \mathbf{y}') \geq \gamma$$

$$\forall (\mathbf{x}_t, \mathbf{y}_t) \in \mathcal{D}$$

and $\mathbf{y}' \in \mathcal{Y}$

Min Norm:

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2$$

such that:

$$\mathbf{w} \cdot \phi(\mathbf{x}_t, \mathbf{y}_t) - \mathbf{w} \cdot \phi(\mathbf{x}_t, \mathbf{y}') \geq 1$$

$$\forall (\mathbf{x}_t, \mathbf{y}_t) \in \mathcal{D}$$

and $\mathbf{y}' \in \mathcal{Y}$

- Instead of fixing $\|\mathbf{w}\|$ we fix the margin $\gamma = 1$
- Make substitution $\mathbf{w}' = \mathbf{w}/\gamma$; then we have $\gamma = \frac{\|\mathbf{w}\|}{\|\mathbf{w}'\|} = \frac{1}{\|\mathbf{w}'\|}$.

Support Vector Machines

$$\mathbf{w} = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2$$

such that:

$$\mathbf{w} \cdot \phi(\mathbf{x}_t, \mathbf{y}_t) - \mathbf{w} \cdot \phi(\mathbf{x}_t, \mathbf{y}') \geq 1$$

$$\forall (\mathbf{x}_t, \mathbf{y}_t) \in \mathcal{D} \text{ and } \mathbf{y}' \in \mathcal{Y}$$

- **Quadratic programming problem** – a well known convex optimization problem
- Can be solved with many techniques.

Support Vector Machines

What if data is not separable?

$$\mathbf{w} = \arg \min_{\mathbf{w}, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{t=1}^N \xi_t$$

such that:

$$\mathbf{w} \cdot \phi(\mathbf{x}_t, \mathbf{y}_t) - \mathbf{w} \cdot \phi(\mathbf{x}_t, \mathbf{y}') \geq 1 - \xi_t \text{ and } \xi_t \geq 0$$

$$\forall (\mathbf{x}_t, \mathbf{y}_t) \in \mathcal{D} \text{ and } \mathbf{y}' \in \mathcal{Y}$$

ξ_t : trade-off between margin per example and $\|\mathbf{w}\|$
Larger C = more examples correctly classified

Kernels

Historically, SVMs with kernels co-occurred together and were extremely popular

Can “kernelize” algorithms to make them non-linear (not only SVMs, but also logistic regression, perceptron, ...)

More later.

Support Vector Machines

$$\mathbf{w} = \arg \min_{\mathbf{w}, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{t=1}^N \xi_t$$

such that:

$$\mathbf{w} \cdot \phi(\mathbf{x}_t, \mathbf{y}_t) - \mathbf{w} \cdot \phi(\mathbf{x}_t, \mathbf{y}') \geq 1 - \xi_t$$

Support Vector Machines

$$\mathbf{w} = \arg \min_{\mathbf{w}, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{t=1}^N \xi_t$$

such that:

$$\mathbf{w} \cdot \phi(\mathbf{x}_t, \mathbf{y}_t) - \max_{\mathbf{y}' \neq \mathbf{y}_t} \mathbf{w} \cdot \phi(\mathbf{x}_t, \mathbf{y}') \geq 1 - \xi_t$$

Support Vector Machines

$$\mathbf{w} = \arg \min_{\mathbf{w}, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{t=1}^N \xi_t$$

such that:

$$\xi_t \geq 1 + \max_{y' \neq y_t} \mathbf{w} \cdot \phi(x_t, y') - \mathbf{w} \cdot \phi(x_t, y_t)$$

Support Vector Machines

$$\mathbf{w} = \arg \min_{\mathbf{w}, \xi} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \sum_{t=1}^N \xi_t \quad \lambda = \frac{1}{C}$$

such that:

$$\xi_t \geq 1 + \max_{y' \neq y_t} \mathbf{w} \cdot \phi(x_t, y') - \mathbf{w} \cdot \phi(x_t, y_t)$$

Support Vector Machines

$$\mathbf{w} = \arg \min_{\mathbf{w}, \xi} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \sum_{t=1}^N \xi_t \quad \lambda = \frac{1}{C}$$

such that:

$$\xi_t \geq 1 + \max_{y' \neq y_t} \mathbf{w} \cdot \phi(\mathbf{x}_t, \mathbf{y}') - \mathbf{w} \cdot \phi(\mathbf{x}_t, \mathbf{y}_t)$$

If $\|\mathbf{w}\|$ classifies $(\mathbf{x}_t, \mathbf{y}_t)$ with margin 1, penalty $\xi_t = 0$

Otherwise penalty $\xi_t = 1 + \max_{y' \neq y_t} \mathbf{w} \cdot \phi(\mathbf{x}_t, \mathbf{y}') - \mathbf{w} \cdot \phi(\mathbf{x}_t, \mathbf{y}_t)$

Support Vector Machines

$$\mathbf{w} = \arg \min_{\mathbf{w}, \xi} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \sum_{t=1}^N \xi_t \quad \lambda = \frac{1}{C}$$

such that:

$$\xi_t \geq 1 + \max_{y' \neq y_t} \mathbf{w} \cdot \phi(x_t, y') - \mathbf{w} \cdot \phi(x_t, y_t)$$

If $\|\mathbf{w}\|$ classifies (x_t, y_t) with margin 1, penalty $\xi_t = 0$

Otherwise penalty $\xi_t = 1 + \max_{y' \neq y_t} \mathbf{w} \cdot \phi(x_t, y') - \mathbf{w} \cdot \phi(x_t, y_t)$

Hinge loss:

$$L((x_t, y_t); \mathbf{w}) = \max(0, 1 + \max_{y' \neq y_t} \mathbf{w} \cdot \phi(x_t, y') - \mathbf{w} \cdot \phi(x_t, y_t))$$

Support Vector Machines

$$\mathbf{w} = \arg \min_{\mathbf{w}, \xi} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \sum_{t=1}^N \xi_t$$

such that:

$$\xi_t \geq 1 + \max_{y' \neq y_t} \mathbf{w} \cdot \phi(x_t, y') - \mathbf{w} \cdot \phi(x_t, y_t)$$

Hinge loss equivalent

$$\mathbf{w} = \arg \min_{\mathbf{w}} \sum_{t=1}^N L((x_t, y_t); \mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

$$= \arg \min_{\mathbf{w}} \left(\sum_{t=1}^N \max (0, 1 + \max_{y' \neq y_t} \mathbf{w} \cdot \phi(x_t, y') - \mathbf{w} \cdot \phi(x_t, y_t)) \right) + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

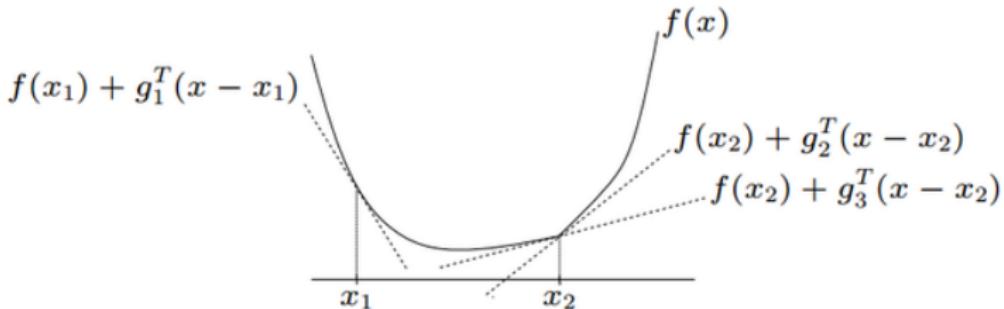
From Gradient to Subgradient

The hinge loss is a **piecewise linear function**—not differentiable everywhere

Cannot use gradient descent

But... can use **subgradient** descent (almost the same)!

Recap: Subgradient



- Defined for convex functions $f : \mathbb{R}^D \rightarrow \mathbb{R}$
- Generalizes the notion of gradient—in points where f is differentiable, there is a single subgradient which equals the gradient
- Other points may have multiple subgradients

Subgradient Descent

$$\begin{aligned} L(\mathbf{w}; (x, y)) &= \max (0, 1 + \max_{y' \neq y} \mathbf{w} \cdot \phi(x, y') - \mathbf{w} \cdot \phi(x, y)) \\ &= \left(\max_{y' \in \mathcal{Y}} \mathbf{w} \cdot \phi(x, y') + [[y' \neq y]] \right) - \mathbf{w} \cdot \phi(x, y) \end{aligned}$$

A **subgradient** of the hinge is

$$\partial_{\mathbf{w}} L(\mathbf{w}; (x, y)) \ni \phi(x, \hat{y}) - \phi(x, y)$$

where

$$\hat{y} = \arg \max_{y' \in \mathcal{Y}} \mathbf{w} \cdot \phi(x, y') + [[y' \neq y]]$$

Can also train SVMs with (stochastic) sub-gradient descent!

Perceptron and Hinge-Loss

SVM subgradient update looks like perceptron update

$$\mathbf{w}^{k+1} = \mathbf{w}^k - \eta \begin{cases} 0, & \text{if } \mathbf{w} \cdot \phi(\mathbf{x}_t, \mathbf{y}_t) - \max_{\mathbf{y} \neq \mathbf{y}_t} \mathbf{w} \cdot \phi(\mathbf{x}_t, \mathbf{y}) \geq 1 \\ \phi(\mathbf{x}_t, \mathbf{y}) - \phi(\mathbf{x}_t, \mathbf{y}_t), & \text{otherwise, where } \mathbf{y} = \arg \max_{\mathbf{y}} \mathbf{w} \cdot \phi(\mathbf{x}_t, \mathbf{y}) + [[\mathbf{y} \neq \mathbf{y}_t]] \end{cases}$$

Perceptron

$$\mathbf{w}^{k+1} = \mathbf{w}^k - \eta \begin{cases} 0, & \text{if } \mathbf{w} \cdot \phi(\mathbf{x}_t, \mathbf{y}_t) - \max_{\mathbf{y}} \mathbf{w} \cdot \phi(\mathbf{x}_t, \mathbf{y}) \geq 0 \\ \phi(\mathbf{x}_t, \mathbf{y}) - \phi(\mathbf{x}_t, \mathbf{y}_t), & \text{otherwise, where } \mathbf{y} = \arg \max_{\mathbf{y}} \mathbf{w} \cdot \phi(\mathbf{x}_t, \mathbf{y}) \end{cases}$$

where $\eta = 1$

Perceptron = SGD with no-margin hinge-loss

$$\max (0, 1 + \max_{\mathbf{y} \neq \mathbf{y}_t} \mathbf{w} \cdot \phi(\mathbf{x}_t, \mathbf{y}) - \mathbf{w} \cdot \phi(\mathbf{x}_t, \mathbf{y}_t))$$

What we have covered

- Linear Classifiers
 - Naive Bayes
 - Logistic Regression
 - Perceptron
 - Support Vector Machines

What is next

- Regularization
- Softmax and sparsemax
- Non-linear classifiers

Outline

① Data and Feature Representation

② Perceptron

③ Naive Bayes

④ Logistic Regression

⑤ Support Vector Machines

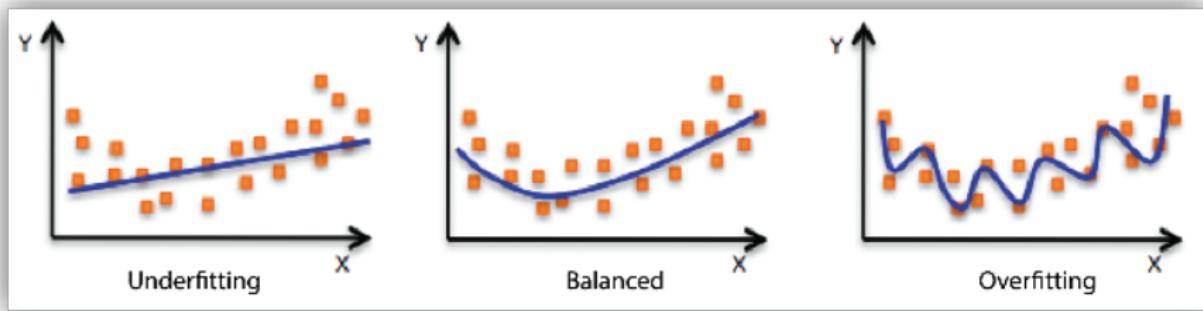
⑥ Regularization

⑦ Non-Linear Classifiers

Regularization

Overfitting

If the model is too complex (too many parameters) and the data is scarce, we run the risk of **overfitting**:



- We saw one example already when talking about add-one smoothing in Naive Bayes!

Regularization

In practice, we **regularize** models to prevent overfitting

$$\arg \min_{\mathbf{w}} \sum_{t=1}^N L(\mathbf{w}; (x_t, y_t)) + \lambda \Omega(\mathbf{w}),$$

where $\Omega(\mathbf{w})$ is the regularization function, and λ controls how much to regularize.

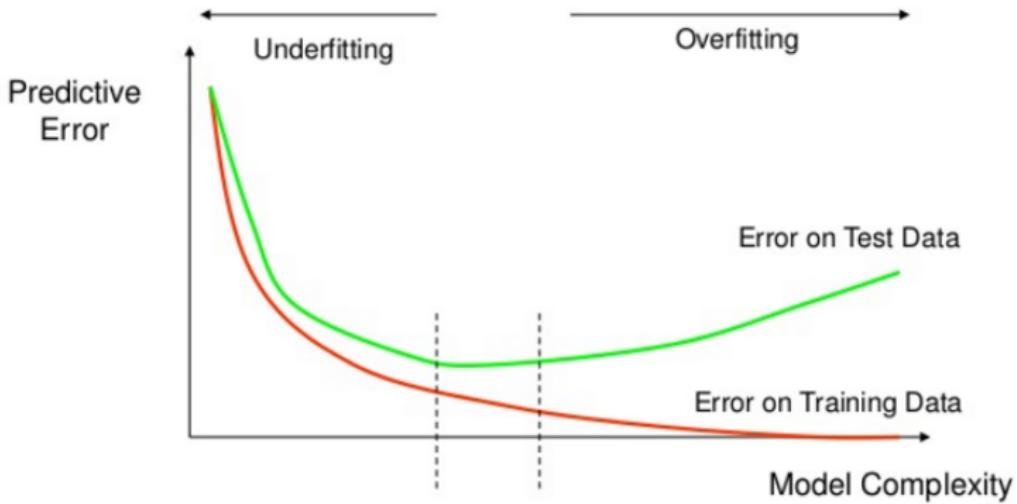
- Gaussian prior (ℓ_2), promotes smaller weights:

$$\Omega(\mathbf{w}) = \|\mathbf{w}\|_2^2 = \sum_i w_i^2.$$

- Laplacian prior (ℓ_1), promotes **sparse** weights!

$$\Omega(\mathbf{w}) = \|\mathbf{w}\|_1 = \sum_i |w_i|$$

Empirical Risk Minimization



Logistic Regression with ℓ_2 Regularization

$$\sum_{t=1}^N L(\mathbf{w}; (\mathbf{x}_t, \mathbf{y}_t)) + \lambda \Omega(\mathbf{w}) = -\sum_{t=1}^N \log(\exp(\mathbf{w} \cdot \phi(\mathbf{x}_t, \mathbf{y}_t)) / Z_x) + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

- What is the new gradient?

$$\sum_{t=1}^N \nabla_{\mathbf{w}} L(\mathbf{w}; (\mathbf{x}_t, \mathbf{y}_t)) + \nabla_{\mathbf{w}} \lambda \Omega(\mathbf{w})$$

- We know $\nabla_{\mathbf{w}} L(\mathbf{w}; (\mathbf{x}_t, \mathbf{y}_t))$
- Just need $\nabla_{\mathbf{w}} \frac{\lambda}{2} \|\mathbf{w}\|^2 = \lambda \mathbf{w}$

Support Vector Machines

Hinge-loss formulation: ℓ_2 regularization already happening!

$$\begin{aligned}\mathbf{w} &= \arg \min_{\mathbf{w}} \sum_{t=1}^N L(\mathbf{w}; (\mathbf{x}_t, \mathbf{y}_t)) + \lambda \Omega(\mathbf{w}) \\ &= \arg \min_{\mathbf{w}} \sum_{t=1}^N \max (0, 1 + \max_{y \neq y_t} \mathbf{w} \cdot \phi(\mathbf{x}_t, y) - \mathbf{w} \cdot \phi(\mathbf{x}_t, y_t)) + \lambda \Omega(\mathbf{w}) \\ &= \arg \min_{\mathbf{w}} \sum_{t=1}^N \max (0, 1 + \max_{y \neq y_t} \mathbf{w} \cdot \phi(\mathbf{x}_t, y) - \mathbf{w} \cdot \phi(\mathbf{x}_t, y_t)) + \frac{\lambda}{2} \|\mathbf{w}\|^2\end{aligned}$$

↑ SVM optimization ↑

SVMs vs. Logistic Regression

$$\mathbf{w} = \arg \min_{\mathbf{w}} \sum_{t=1}^N L(\mathbf{w}; (\mathbf{x}_t, y_t)) + \lambda \Omega(\mathbf{w})$$

SVMs vs. Logistic Regression

$$\mathbf{w} = \arg \min_{\mathbf{w}} \sum_{t=1}^N L(\mathbf{w}; (\mathbf{x}_t, y_t)) + \lambda \Omega(\mathbf{w})$$

SVMs/hinge-loss: $\max (0, 1 + \max_{y \neq y_t} (\mathbf{w} \cdot \phi(\mathbf{x}_t, y) - \mathbf{w} \cdot \phi(\mathbf{x}_t, y_t)))$

$$\mathbf{w} = \arg \min_{\mathbf{w}} \sum_{t=1}^N \max (0, 1 + \max_{y \neq y_t} \mathbf{w} \cdot \phi(\mathbf{x}_t, y) - \mathbf{w} \cdot \phi(\mathbf{x}_t, y_t)) + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

SVMs vs. Logistic Regression

$$\mathbf{w} = \arg \min_{\mathbf{w}} \sum_{t=1}^N L(\mathbf{w}; (\mathbf{x}_t, y_t)) + \lambda \Omega(\mathbf{w})$$

SVMs/hinge-loss: $\max (0, 1 + \max_{y \neq y_t} (\mathbf{w} \cdot \phi(\mathbf{x}_t, y) - \mathbf{w} \cdot \phi(\mathbf{x}_t, y_t)))$

$$\mathbf{w} = \arg \min_{\mathbf{w}} \sum_{t=1}^N \max (0, 1 + \max_{y \neq y_t} \mathbf{w} \cdot \phi(\mathbf{x}_t, y) - \mathbf{w} \cdot \phi(\mathbf{x}_t, y_t)) + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

Logistic Regression/**log-loss**: $-\log (\exp(\mathbf{w} \cdot \phi(\mathbf{x}_t, y_t))/Z_x)$

$$\mathbf{w} = \arg \min_{\mathbf{w}} \sum_{t=1}^N -\log (\exp(\mathbf{w} \cdot \phi(\mathbf{x}_t, y_t))/Z_x) + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

Loss Function

Should match as much as possible the metric we want to optimize at test time

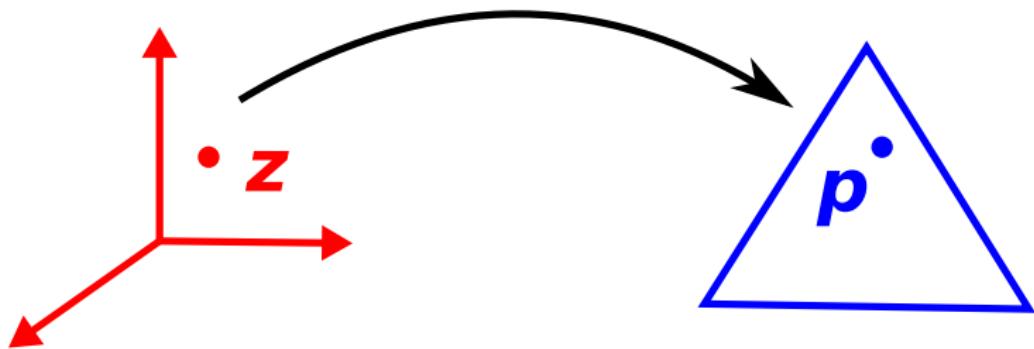
Should be well-behaved (continuous, maybe smooth) to be amenable to optimization (this rules out the 0/1 loss)

Some examples:

- Squared loss for regression
- Negative log-likelihood (cross-entropy): multinomial logistic regression
- Hinge loss: support vector machines
- Sparsemax loss for multi-class and multi-label classification ([next](#))

Recap

How to map from a set of label scores $\mathbb{R}^{|\mathcal{Y}|}$ to a probability distribution over \mathcal{Y} ?



We already saw one example: softmax.

Next: sparsemax.

Recap: Softmax Transformation

The typical transformation for multi-class classification is
softmax : $\mathbb{R}^{|\mathcal{Y}|} \rightarrow \Delta^{|\mathcal{Y}|-1}$:

$$\text{softmax}(z) = \left[\frac{\exp(z_1)}{\sum_c \exp(z_c)}, \dots, \frac{\exp(z_{|\mathcal{Y}|})}{\sum_c \exp(z_c)} \right]$$

- Underlies multinomial logistic regression!
- Strictly positive, sums to 1
- Resulting distribution has full support: $\text{softmax}(z) > \mathbf{0}, \forall z$
- A disadvantage if a *sparse* probability distribution is desired
- Common workaround: threshold and truncate

Sparsemax (Martins and Astudillo, 2016)

A sparse-friendly alternative is **sparsemax** : $\mathbb{R}^{|\mathcal{Y}|} \rightarrow \Delta^{|\mathcal{Y}|-1}$, defined as:

$$\text{sparsemax}(z) := \arg \min_{\mathbf{p} \in \Delta^{|\mathcal{Y}|-1}} \|\mathbf{p} - z\|^2.$$

- In words: Euclidean projection of z onto the probability simplex
- Likely to hit the boundary of the simplex, in which case $\text{sparsemax}(z)$ becomes sparse (hence the name)
- Retains many of the properties of softmax (e.g. differentiability), having in addition the ability of producing sparse distributions
- Projecting onto the simplex amounts to a **soft-thresholding** operation
- Efficient linear time forward/backward propagation (see paper)

Sparsemax in Closed Form

- Projecting onto the simplex amounts to a soft-thresholding operation:

$$\text{sparsemax}_i(\mathbf{z}) = \max\{0, z_i - \tau\}$$

where τ is a normalizing constant such that $\sum_j \max\{0, z_j - \tau\} = 1$

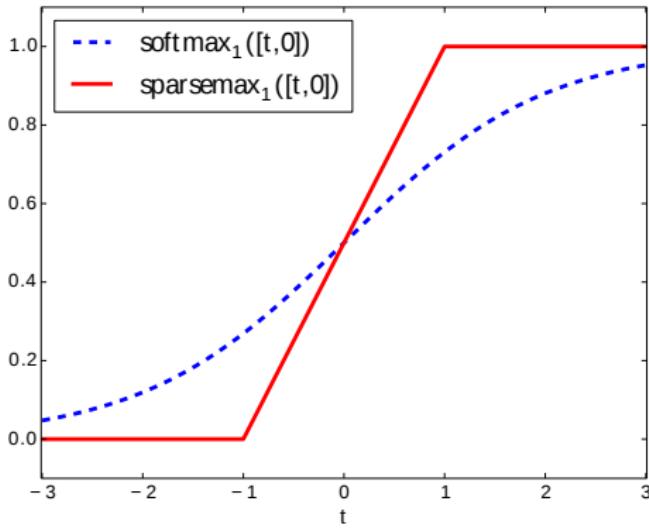
- To evaluate the sparsemax, all we need is to compute τ
- Coordinates above the threshold will be shifted by this amount; the others will be truncated to zero

Two Dimensions

- Parametrize $z = (t, 0)$
- The 2D **softmax** is the logistic (sigmoid) function:

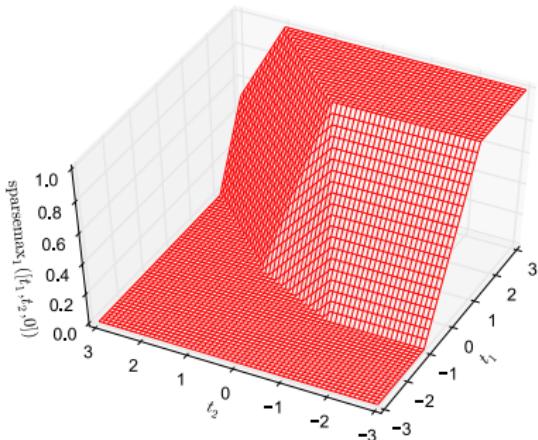
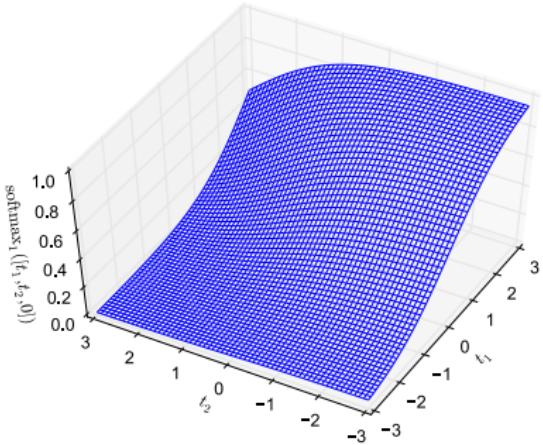
$$\text{softmax}_1(z) = (1 + \exp(-t))^{-1}$$

- The 2D **sparsemax** is the “hard” version of the sigmoid:



Three Dimensions

- Parameterize $z = (t_1, t_2, 0)$ and plot **softmax**₁(z) and **sparsemax**₁(z) as a function of t_1 and t_2
- **sparsemax** is piecewise linear, but asymptotically similar to **softmax**



Loss Function

How to use sparsemax as a loss function?

Caveat: sparsemax is sparse and we don't want to take the log of zero...

Recap: Multinomial Logistic Regression

- The common choice for a softmax output layer
- The classifier estimates $P(y = c \mid \mathbf{x}; \mathbf{w})$
- We minimize the negative log-likelihood:

$$\begin{aligned} L(\mathbf{w}; (\mathbf{x}, \mathbf{y})) &= -\log P(y \mid \mathbf{x}; \mathbf{w}) \\ &= -\log \text{softmax}(z(\mathbf{x})), \end{aligned}$$

where $z_c(\mathbf{x}) = \mathbf{w} \cdot \phi(\mathbf{x}, c)$.

- Loss gradient:

$$\nabla_{\mathbf{w}} L((\mathbf{x}, \mathbf{y}); \mathbf{w}) = - \left(\phi(\mathbf{x}, \mathbf{y}) - \sum_c \text{softmax}_c(z(\mathbf{x})) \phi(\mathbf{x}, c) \right)$$

Sparsemax Loss (Martins and Astudillo, 2016)

- The natural choice for a sparsemax output layer
- The neural network estimates $P(y = c | x; \mathbf{w})$ as a **sparse distribution**

$$L((\mathbf{x}, \mathbf{y}); \mathbf{w}) = -z_c + \frac{1}{2} \sum_{j \in S(z)} (z_j^2 - \tau^2(z)) + \frac{1}{2},$$

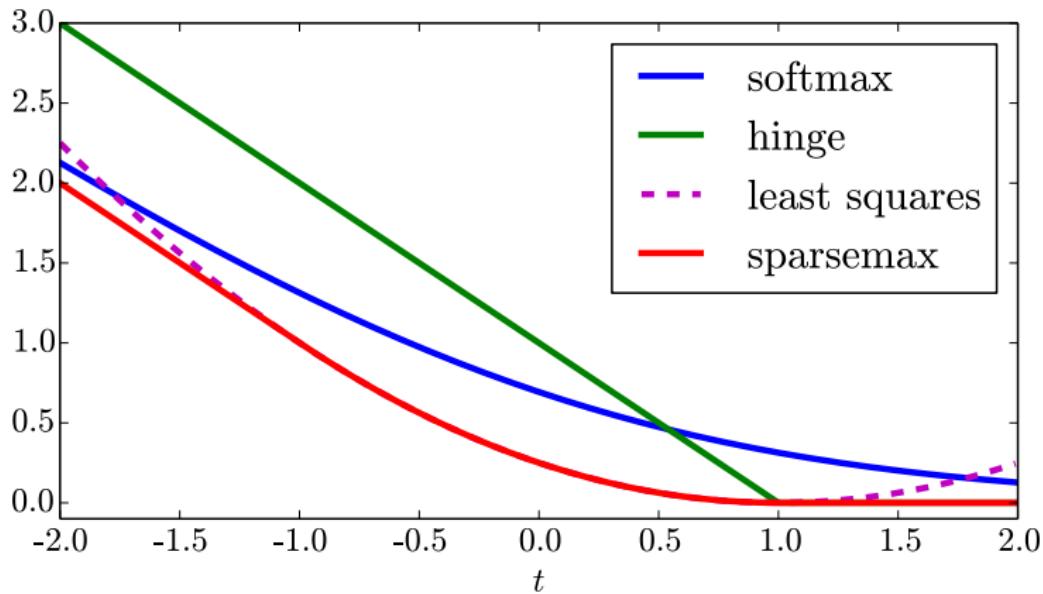
where $z_c(x) = \mathbf{w} \cdot \phi(x, c)$, $S(z)$ is the support of $P(\mathbf{y} | \mathbf{x}; \mathbf{w})$ and $\tau^2 : \mathbb{R}^K \rightarrow \mathbb{R}$ is the square of the threshold function.

- Loss gradient:

$$\nabla_{\mathbf{w}} L((\mathbf{x}, \mathbf{y}); \mathbf{w}) = - \left(\phi(\mathbf{x}, \mathbf{y}) - \sum_c \text{sparsemax}_c(z(\mathbf{x})) \phi(\mathbf{x}, c) \right)$$

Classification Losses in Two Dimensions

- Let the correct label be $y = 1$ and define $t = z_1 - z_2$:



Outline

① Data and Feature Representation

② Perceptron

③ Naive Bayes

④ Logistic Regression

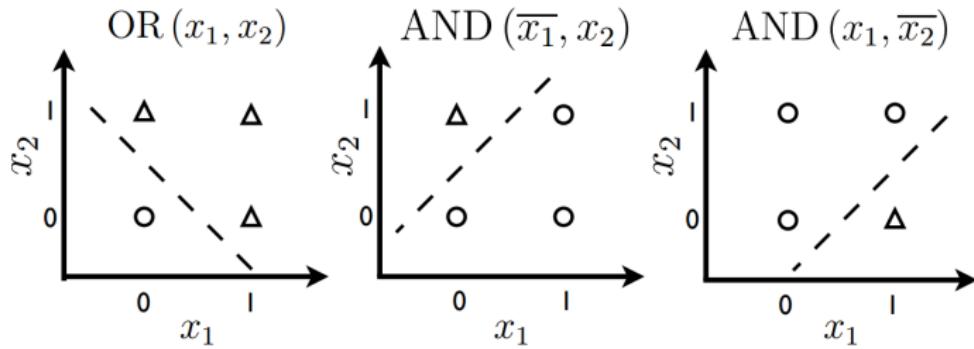
⑤ Support Vector Machines

⑥ Regularization

⑦ Non-Linear Classifiers

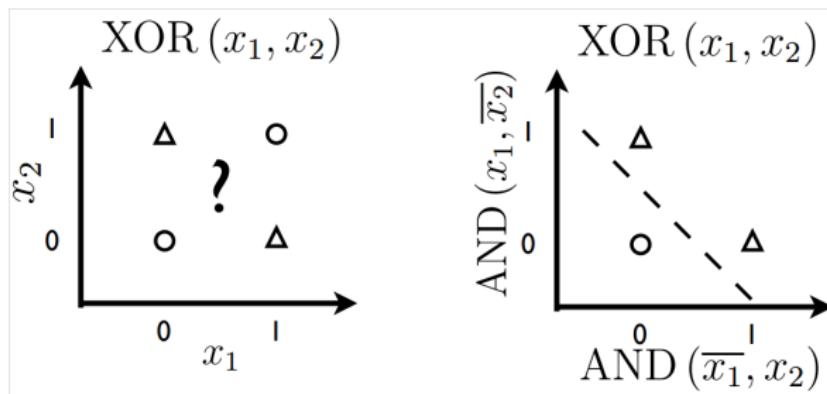
Recap: What a Linear Classifier Can Do

- It **can** solve linearly separable problems (OR, AND)



Recap: What a Linear Classifier **Can't** Do

- ... but it **can't** solve **non-linearly separable** problems such as simple XOR (unless input is transformed into a better representation):



- This was observed by Minsky and Papert (1969) (for the perceptron) and motivated strong criticisms

Summary: Linear Classifiers

We've seen

- Perceptron
- Naive Bayes
- Logistic regression
- Support vector machines

All lead to **convex** optimization problems \Rightarrow no issues with local minima/initialization

All assume the features are well-engineered such that **the data is nearly linearly separable**

What If Data Are Not Linearly Separable?

What If Data Are Not Linearly Separable?

Engineer better features (often works!)



What If Data Are Not Linearly Separable?

Engineer better features (often works!)



Kernel methods:

- works implicitly in a high-dimensional feature space
- ... but still need to choose/design a good kernel
- model capacity confined to positive-definite kernels



What If Data Are Not Linearly Separable?

Engineer better features (often works!)



Kernel methods:

- works implicitly in a high-dimensional feature space
- ... but still need to choose/design a good kernel
- model capacity confined to positive-definite kernels



Neural networks (**next class!**)

- embrace non-convexity and local minima
- instead of engineering features/kernels, engineer the model architecture

Two Views of Machine Learning

There's two big ways of building machine learning systems:

- ① **Feature-based**: describe objects' properties (features) and build models that manipulate them
 - everything that we have seen so far.
- ② **Similarity-based**: don't describe objects by their properties; rather, build systems based on **comparing** objects to each other
 - k -th nearest neighbors; kernel methods; Gaussian processes.

Sometimes the two are equivalent!

Nearest Neighbor Classifier

- Not a linear classifier!
- In its simplest version, doesn't require any parameters
- Instead of "training", **memorize** all the data $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N\}$
- Given a new input x , find its **most similar** data point x_i and predict

$$\hat{y} = y_i$$

- Many variants (e.g. k -th nearest neighbor)
- **Disadvantage:** requires searching over the entire training data
- Specialized data structures can be used to speed up search.

Kernels

- A kernel is a similarity function between two points that is **symmetric** and **positive semi-definite**, which we denote by:

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) \in \mathbb{R}$$

- Given dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N\}$, the **Gram matrix** \mathbf{K} is the $N \times N$ matrix defined as:

$$K_{i,j} = \kappa(\mathbf{x}_i, \mathbf{x}_j)$$

- **Symmetric:**

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \kappa(\mathbf{x}_j, \mathbf{x}_i)$$

- **Positive definite:** for all non-zero \mathbf{v}

$$\mathbf{v}\mathbf{K}\mathbf{v}^T \geq 0$$

Kernels

- **Mercer's Theorem:** for any kernel $\kappa : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, there exists some feature mapping $\psi : \mathcal{X} \rightarrow \mathbb{R}^{\mathcal{X}}$, s.t.:

$$\kappa(x_i, x_j) = \psi(x_i) \cdot \psi(x_j)$$

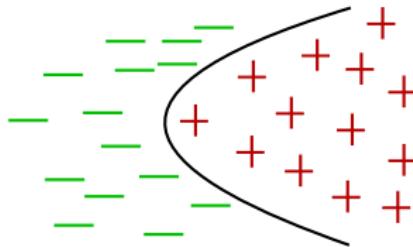
- That is: a kernel corresponds to some a mapping in some **implicit** feature space!
- **Kernel trick:** take a feature-based algorithm (SVMs, perceptron, logistic regression) and replace all explicit feature computations by **kernel evaluations!**

$$w_y \cdot \psi(x) = \sum_{i=1}^N \sum_{y \in \mathcal{Y}} \alpha_{i,y} \kappa(x, x_i) \quad \text{for some } \alpha_{i,y} \in \mathbb{R}$$

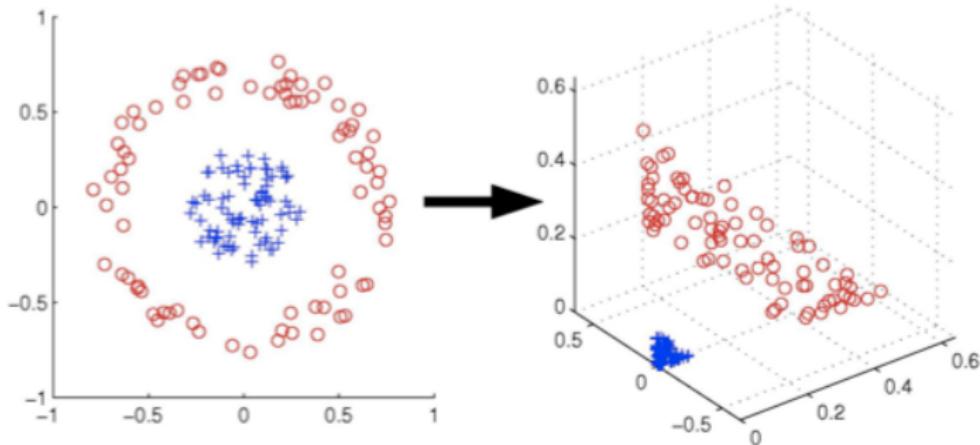
- Extremely popular idea in the 1990-2000s!

Kernels = Tractable Non-Linearity

- A linear classifier in a higher dimensional feature space is a non-linear classifier in the original space
- Computing a non-linear kernel is sometimes better computationally than calculating the corresponding dot product in the high dimension feature space
- Many models can be “kernelized” – learning algorithms generally solve the **dual** optimization problem (also convex)
- Drawback: **quadratic** dependency on dataset size



Linear Classifiers in High Dimension



$$\mathbb{R}^2 \longrightarrow \mathbb{R}^3$$

$$(x_1, x_2) \longmapsto (z_1, z_2, z_3) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$$

Popular Kernels

- Polynomial kernel

$$\kappa(x_i, x_j) = (\psi(x_i) \cdot \psi(x_j) + 1)^d$$

- Gaussian radial basis kernel

$$\kappa(x_i, x_j) = \exp\left(\frac{-||\psi(x_i) - \psi(x_j)||^2}{2\sigma}\right)$$

- String kernels (Lodhi et al., 2002; Collins and Duffy, 2002)
- Tree kernels (Collins and Duffy, 2002)

Conclusions

- Linear classifiers are a broad class including well-known ML methods such as **perceptron**, **Naive Bayes**, **logistic regression**, **support vector machines**
- They all involve manipulating weights and features
- They either lead to closed-form solutions or **convex** optimization problems (**no local minima**)
- Stochastic gradient descent algorithms are useful if training datasets are large
- However, they require manual specification of feature representations
- **Tomorrow:** methods that are able to **learn internal representations**

Thank you!



References I

- Collins, M. and Duffy, N. (2002). Convolution kernels for natural language. *Advances in Neural Information Processing Systems*, 1:625–632.
- Lodhi, H., Saunders, C., Shawe-Taylor, J., Cristianini, N., and Watkins, C. (2002). Text classification using string kernels. *Journal of Machine Learning Research*, 2:419–444.
- Martins, A. F. T. and Astudillo, R. (2016). From Softmax to Sparsemax: A Sparse Model of Attention and Multi-Label Classification. In *Proc. of the International Conference on Machine Learning*.
- Minsky, M. and Papert, S. (1969). Perceptrons.
- Novikoff, A. B. (1962). On convergence proofs for perceptrons. In *Symposium on the Mathematical Theory of Automata*.
- Rosenblatt, F. (1958). The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological review*, 65(6):386.