

Aprendizaje Supervisado

...

Franco Luque - Matías Marenchino

RoadMap

- **Short review of topics from previous course**
Classification / Regression / Perceptron
- **What supervised learning is**
- **SVMs**
- **Ensemble methods**
Bagging: Random Forests / Boosting
- **Neural Networks**

Motivation

Example 1

- **A credit card company receives applications for new credit cards. Each one has information about an applicant:**
 - salary
 - age
 - marital status
 - Veraz
 - Credit report from BCRA
 - ...
- **Problem:** determine if an application should be approved or rejected

Example 2

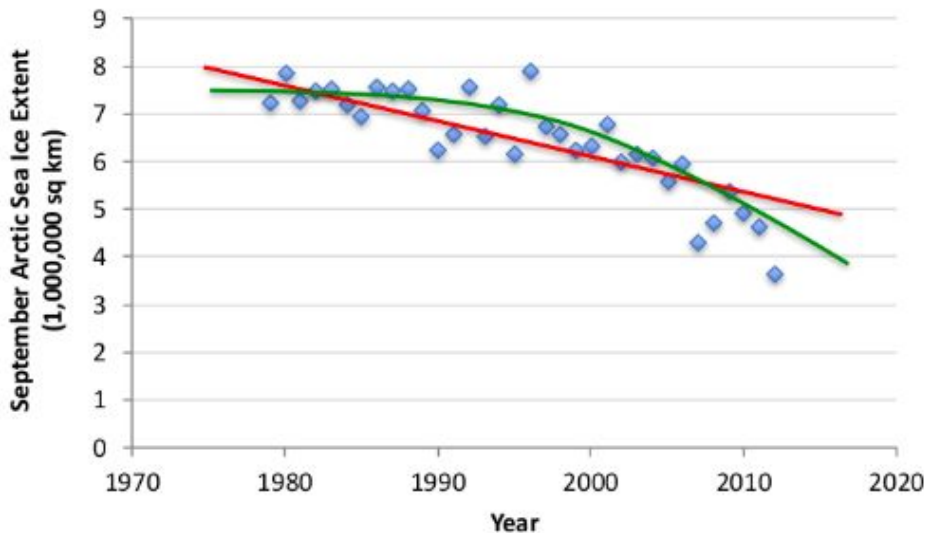
- **Problem:** classify an email as SPAM or not

Describing the problem

- **Data:** A set of records (or samples, instances) described by n attributes: A_1, A_2, \dots, A_n and each sample is labelled with a class (Like SPAM or NOT) or a "score" (like the credit score)
- **Goal:** To learn a model (or a function) from the data that can be used to predict the labels that the records have (and labels for new unlabelled records)

Aprendizaje supervisado: regresión

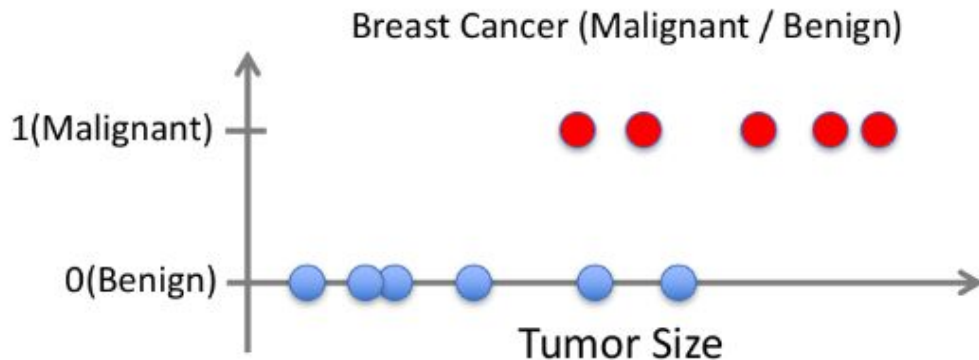
- Datos $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$
- Aprender una $f(x)$ que permita predecir y a partir de x
 - Si y está en $\mathbb{R}^n \rightarrow$ **regresión**



Slides from the previous course

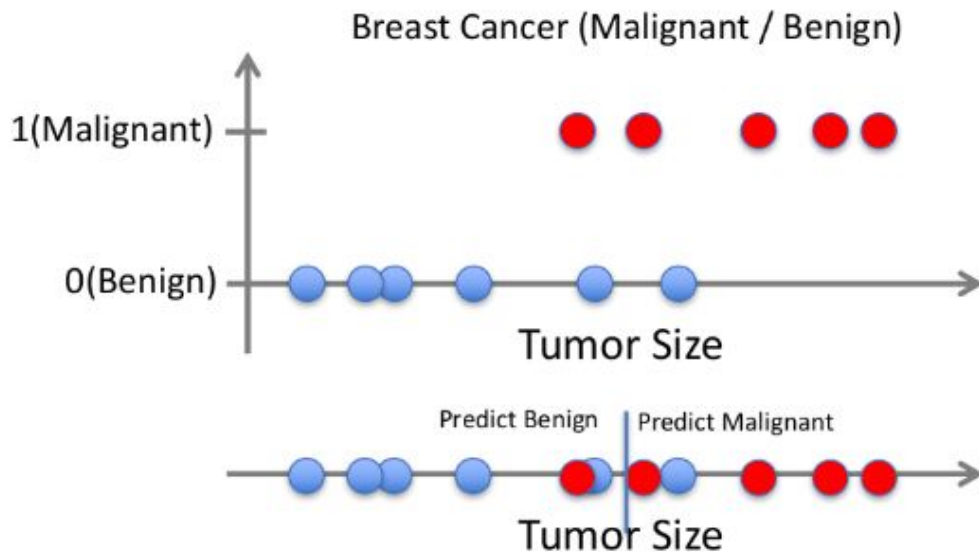
Aprendizaje supervisado: clasificación

- Datos $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$
- Aprender una $f(x)$ que permita predecir y a partir de x
 - Si y es categórica → **clasificación**



Aprendizaje supervisado: clasificación

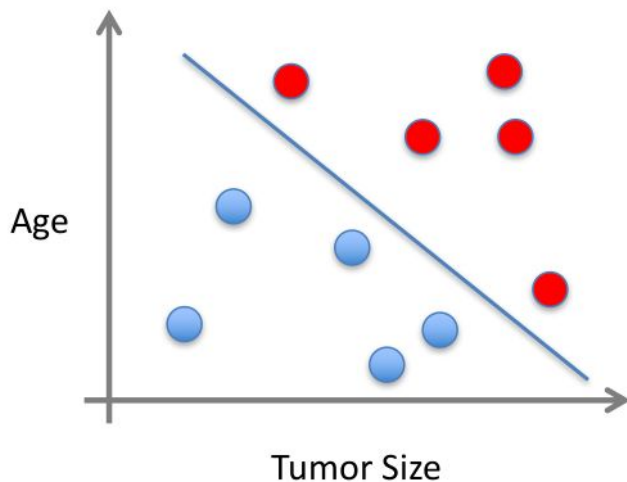
- Datos $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$
- Aprender una $f(x)$ que permita predecir y a partir de x
 - Si y es categórica → **clasificación**



Slides from the previous
course

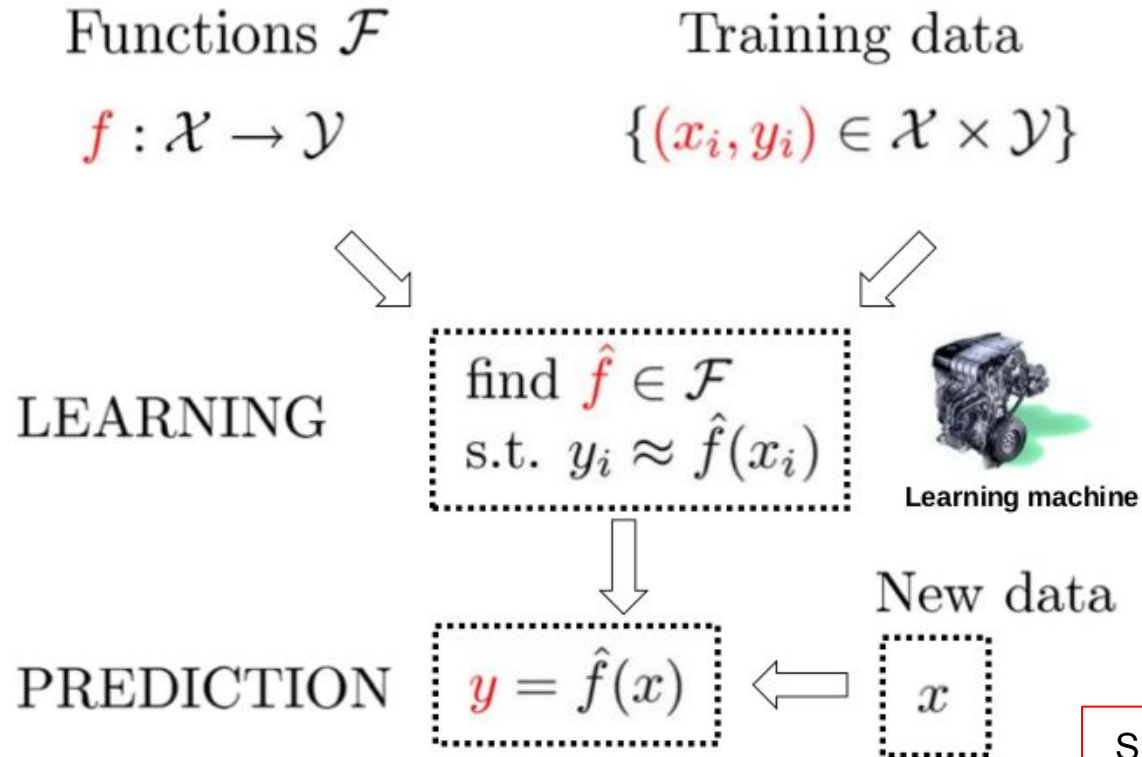
Supervised Learning

- x can be multi-dimensional
 - Each dimension corresponds to an attribute



- Clump Thickness
- Uniformity of Cell Size
- Uniformity of Cell Shape
- ...

Aprendizaje supervisado

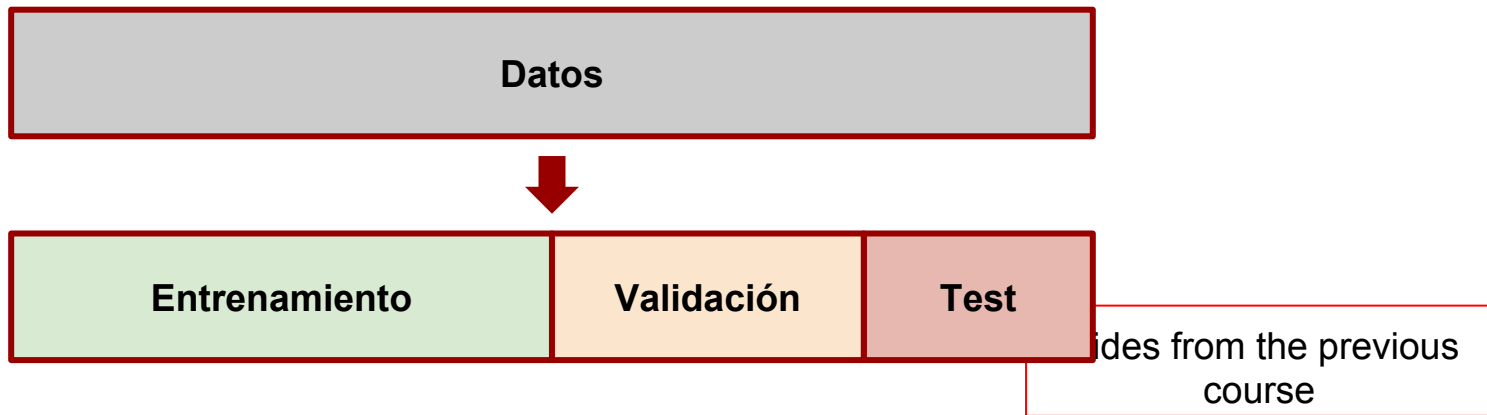


Slides from the previous
course

Elección de *hiperparámetros*

Dividir el conjunto total de ejemplos en tres subconjuntos

- **Entrenamiento:** aprendizaje de variables del modelo
- **Validación:** ajuste/elección de hiperparámetros
- **Test:** estimación final de la performance del modelo entrenado (y con hiperparámetros elegidos adecuadamente)



El algoritmo del "perceptrón"

- Propuesto por Roseblatt en 1958
- El objetivo es encontrar un hiperplano de separación
 - Si los datos son linealmente separables, lo encuentra
- Es un algoritmo *online* (procesa un ejemplo a la vez)
- Muchas variantes ...

El algoritmo del "perceptrón"

Entrada:

- una secuencia de pares de entrenamiento $(x_1, y_1), (x_2, y_2) \dots$
- Una tasa de aprendizaje r

Algoritmo:

- Inicializar $w^{(0)} \in \mathbb{R}^n$
- Para cada ejemplo (x_i, y_i)
 - Predecir $y_i' = \text{sign}(w^T x_i + w_0)$
 - Si $y_i' \neq y_i$:
$$w^{(t+1)} \leftarrow w^{(t)} + r (y_i x_i)$$

El algoritmo del "perceptrón"

Entrada:

- una secuencia de pares de entrenamiento $(x_1, y_1), (x_2, y_2) \dots$
- Una tasa de aprendizaje r (número pequeño y menor a 1)

Algoritmo:

- Inicializar $w^{(0)} \in \mathbb{R}^n$
- Para cada ejemplo (x_i, y_i)
 - Predecir $y_i' = \text{sign}(w^T x_i)$
 - Si $y_i' \neq y_i$:
 $w^{(t+1)} \leftarrow w^{(t)} + r (y_i x_i)$

Actualiza solo cuando comete un error

Error en positivos:

$$w^{(t+1)} \leftarrow w^{(t)} + r x_i$$

Error en negativos:

$$w^{(t+1)} \leftarrow w^{(t)} - r x_i$$

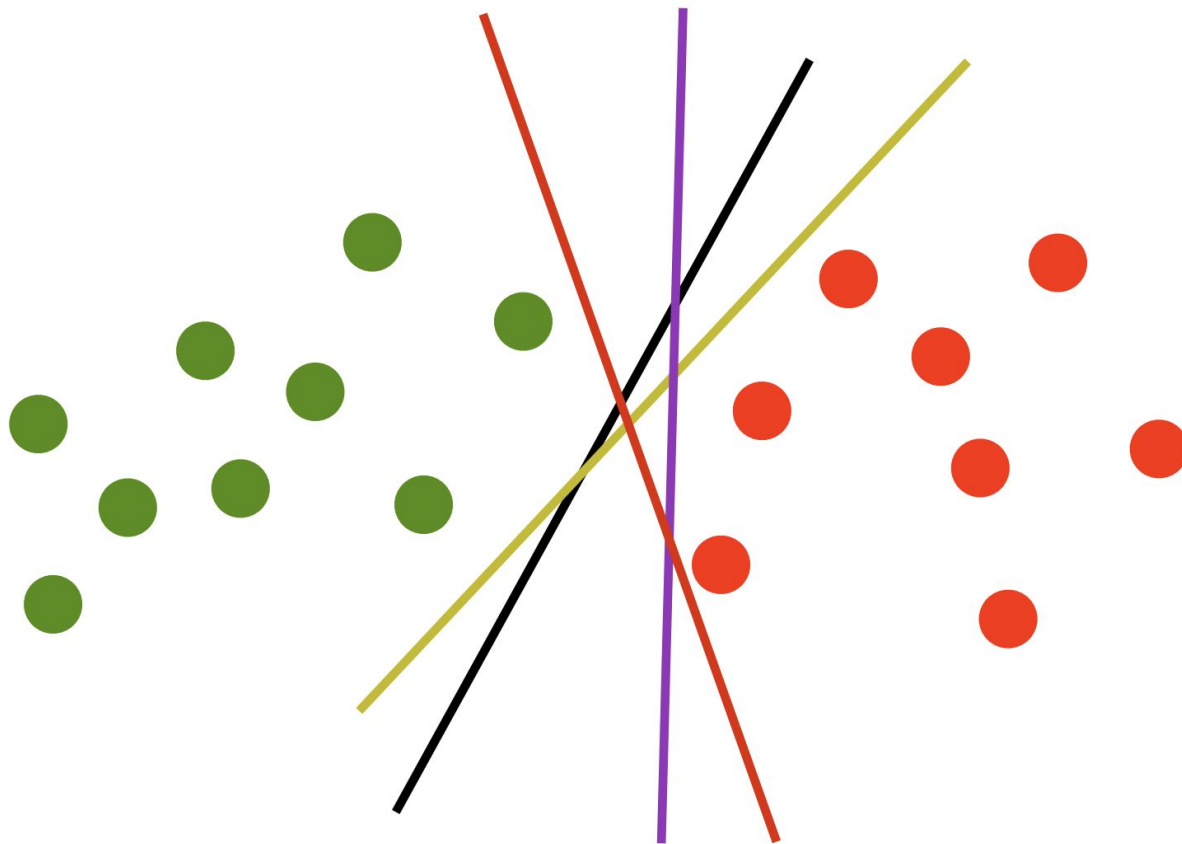
Si $y_i w^T x_i \leq 0 \rightarrow \text{error}$

Slides from the previous
course

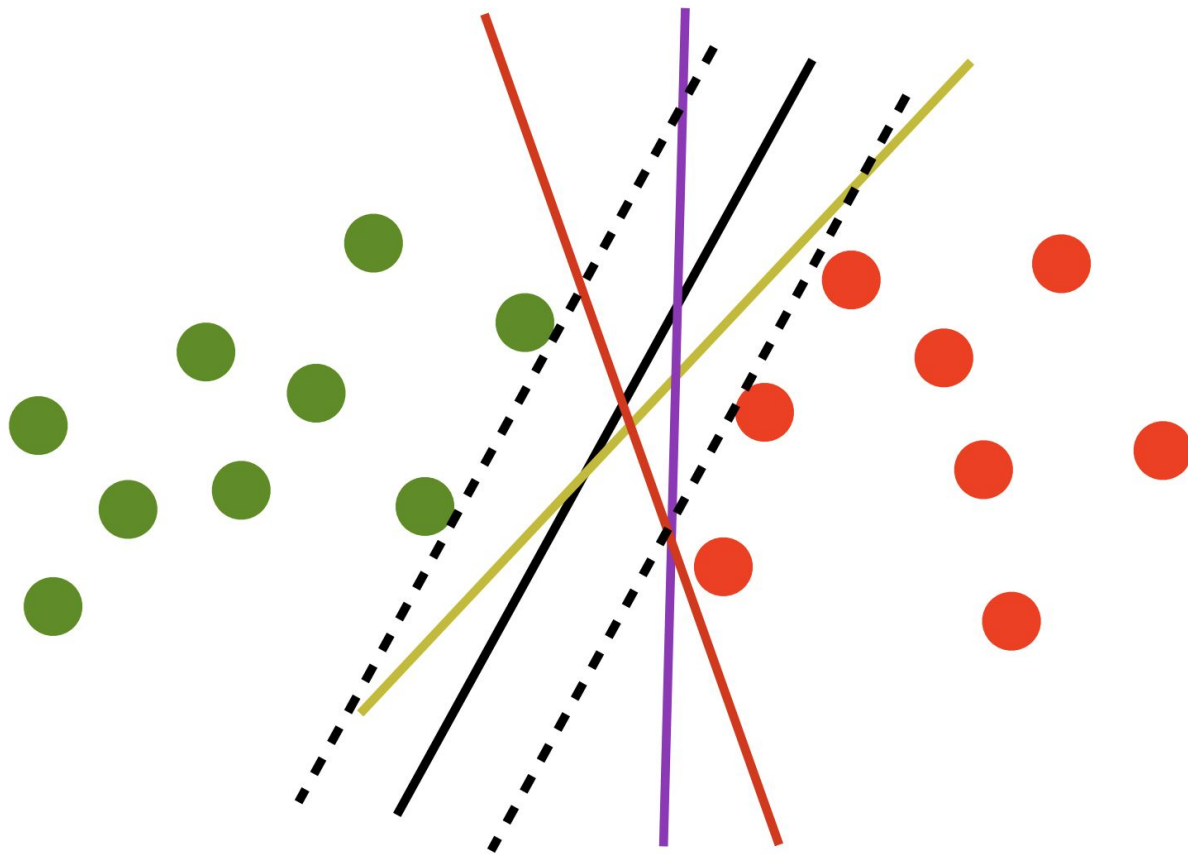
Demo Time (demo 1)

...

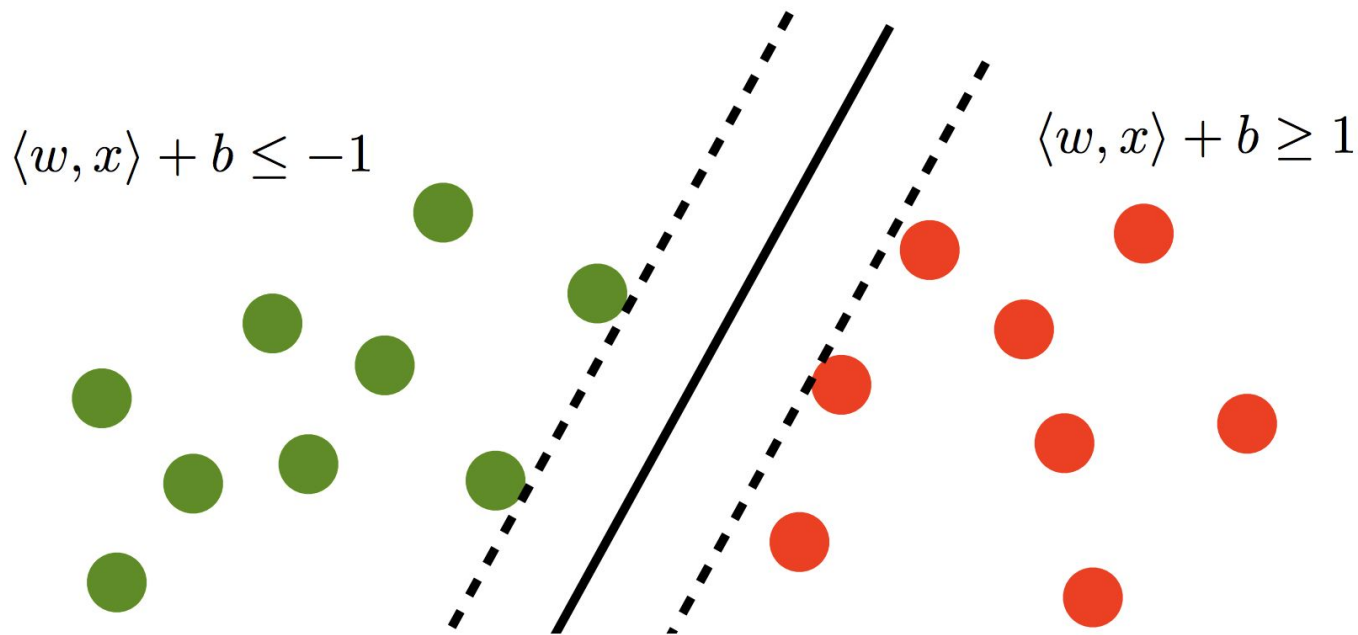
Support Vector Machines



Support Vector Machines



Support Vector Machines



linear function

$$f(x) = \langle w, x \rangle + b$$

Whiteboard Time

...

Support Vector Machines, Deriving the Equations

Support Vector Machines

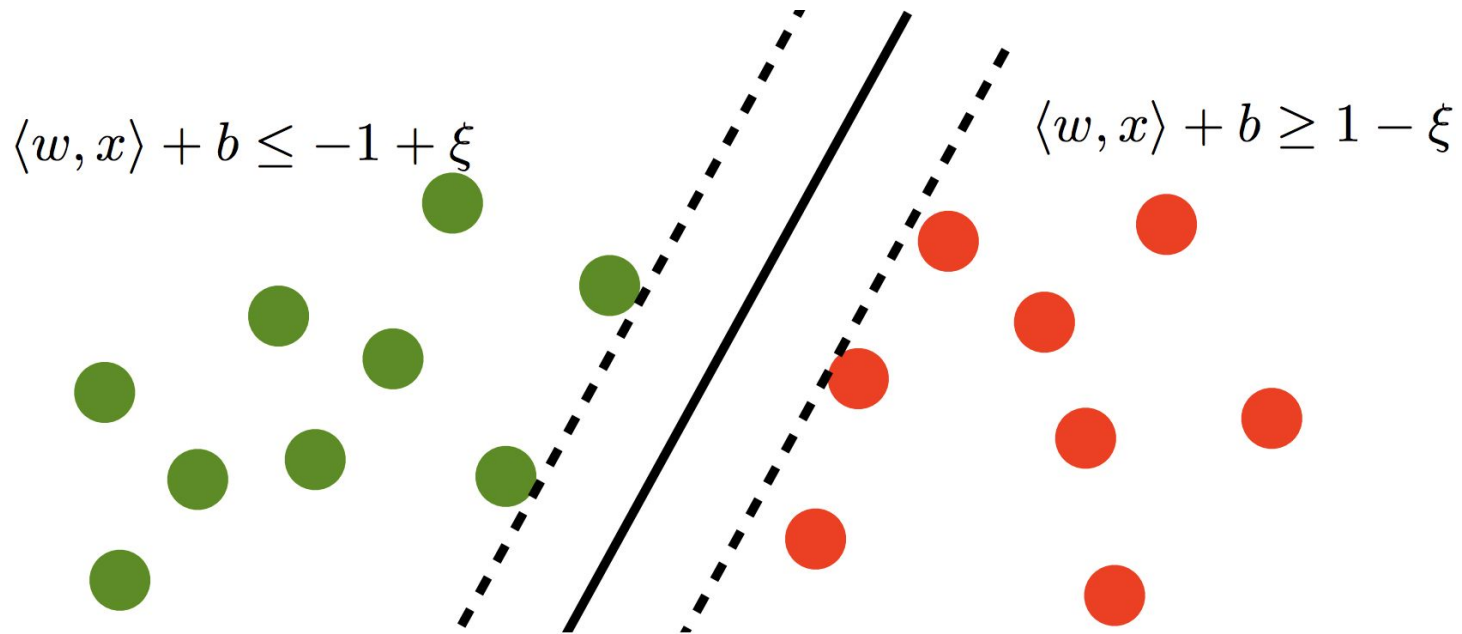
Some Interesting Properties we have found:

- $w = \sum_i \alpha_i y_i x_i$
- Just the points on the margin used to define the hyperplane

Demo Time (demo 2)

...

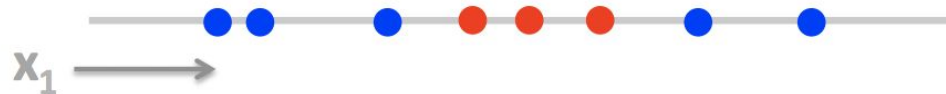
SVMs: slack variables



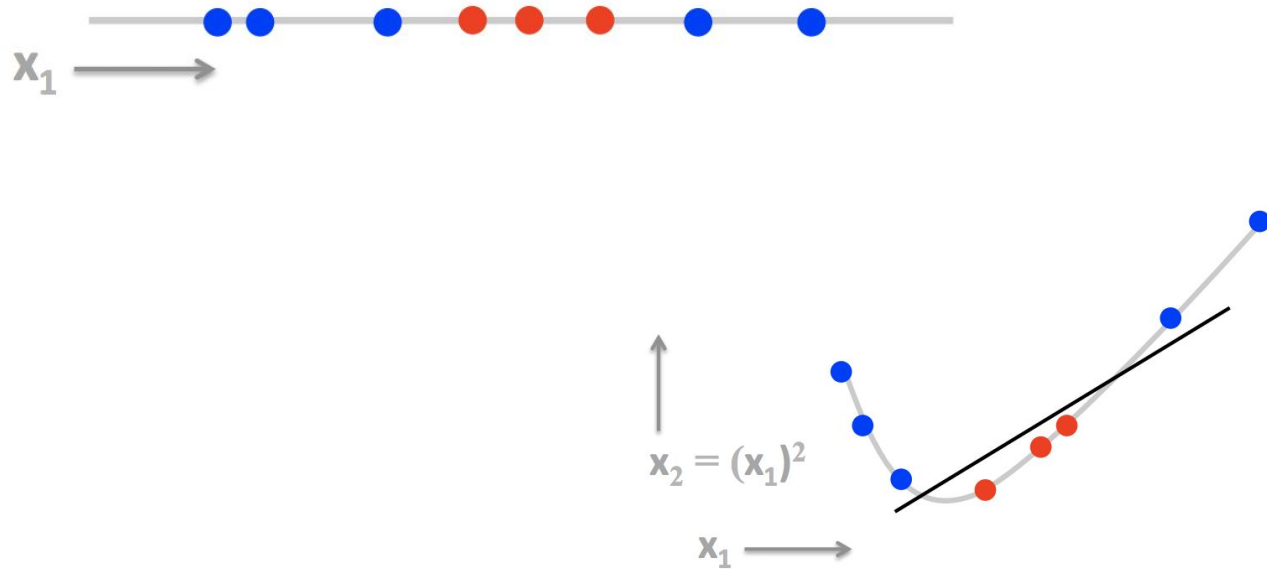
Demo Time (demo 3)

...

SVMs: Kernels



SVMs: Kernels



Whiteboard Time

...

SVMs: Kernels

Some common kernels:

- Polynomial

$$K(x, z) = (1 + \sum_j x_j z_j)^d$$

- Radial Basis Functions (RBF aka Gaussian Similarity Functions)

$$K(x, z) = \exp(-(x - z)^2 / 2\sigma^2)$$

- Sigmoid

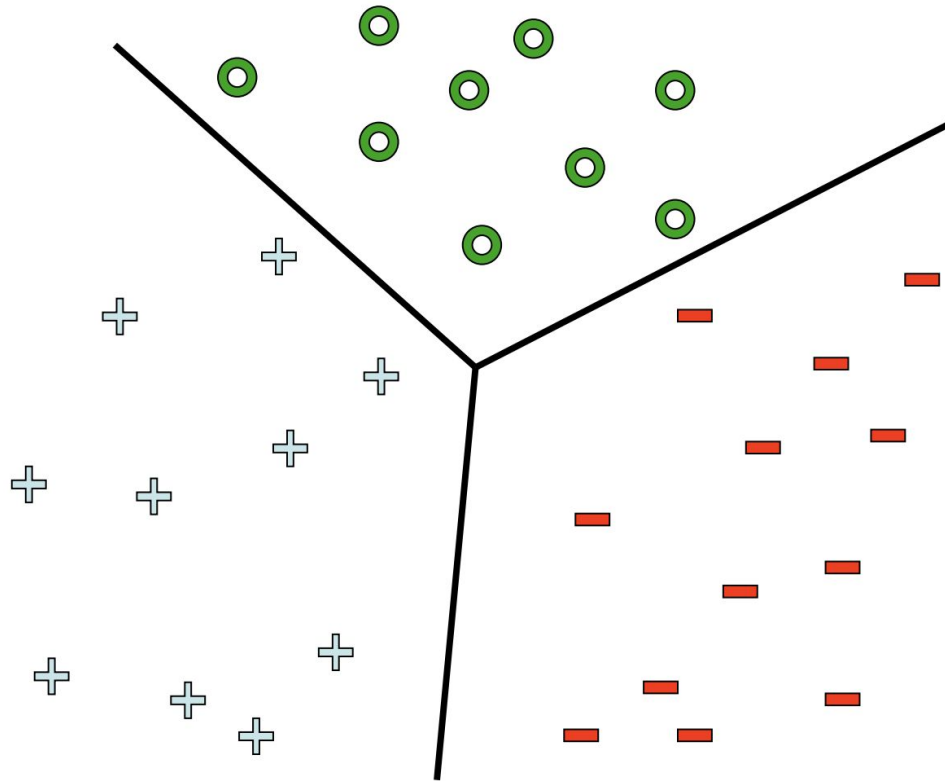
$$K(x, z) = \tanh(c < x, z > + h)$$

- Specific for certain types of problems

Demo Time (demo 4)

...

Multiclass SVMs



Multiclass SVMs: one vs the rest (one vs all)

- **Training:** how could we do it?

Multiclass SVMs: one vs the rest (one vs all)

- **Training:** For M classes:
construct a hyperplane between class k and the other $M - 1$ classes $\Rightarrow M$ SVMs
- **Classification:** how could we do it?

Multiclass SVMs: one vs the rest

- **Training:** For M classes:
construct a hyperplane between class k and the other $M - 1$ classes $\Rightarrow M$ SVMs
- **Classification:** make M predictions (one for each SVMs) and find out the one getting more hits into its positive region.

Multiclass SVMs: one vs one (all vs all)

Multiclass SVMs: one vs one (all vs all)

- **Training:** For M classes:
construct a hyperplane between class i and class j . $\Rightarrow M(M-1)/2$
SVMs
- **Classification:** If f_{ij} is the classifier where i are the positive samples and j the negative ones,, the classification of x is given by

$$f(x) = \operatorname{argmin}_i (\sum_j f_{ij}(x))$$

Support Vector Regressions

- Idea based on support vectors like in SVMs, but now y_i is a real number.
- Uses soft margins in the regression process instead of classification
- Additional parameter ϵ , to compute the loss function

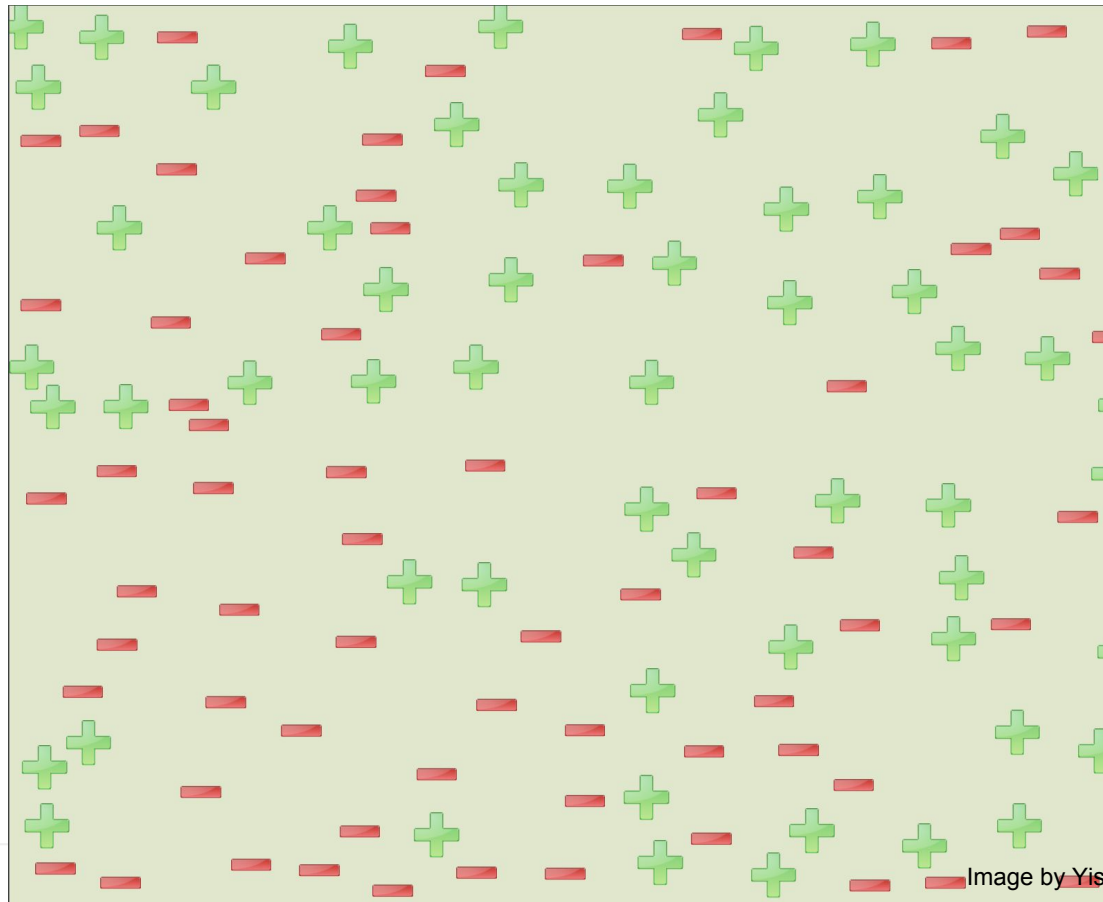
Support Vector Regressions

- Idea based on support vectors like in SVMs, but now y_i is a real number.
- Uses soft margins in the regression process instead of classification
- Additional parameter ε , to compute the loss function
- Not frequently used, logistic regression is more popular

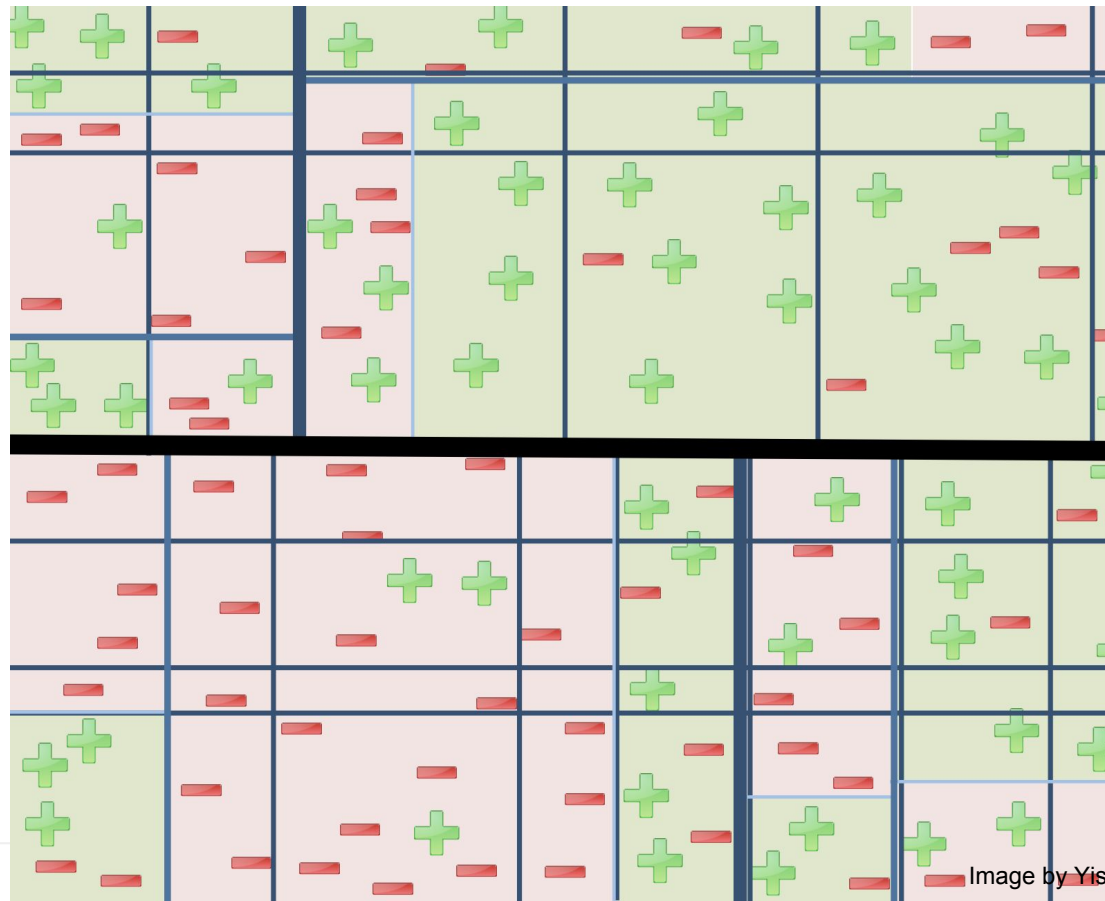
Whiteboard Time

...

Decision Trees (review)



Decision Trees (review)



Decision Trees (review)

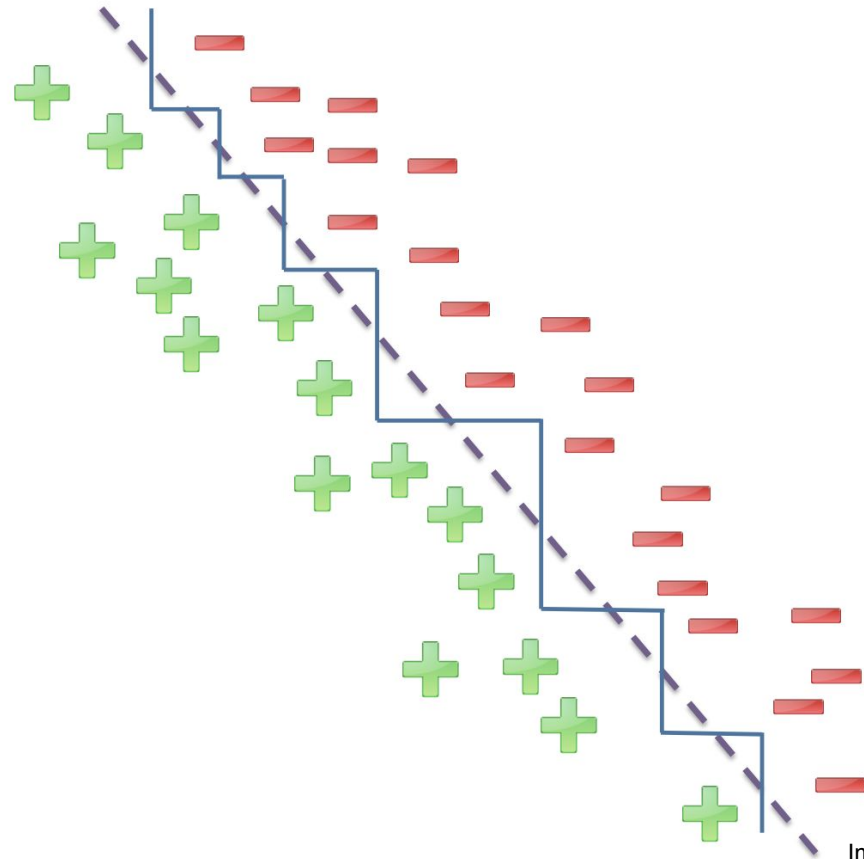


Image by Yisong Yue

Demo Time (demo 5)

...

Ensemble Learning

- **Generate a set of "learners" that, when combined, have higher accuracy.**
- **Assuming we have three learners: L1, L2, L3**
- **The predictions from them may differ**
- **What would we do? Who do we trust?**
 - **Believe the model that we know is best?**
 - **Go with the majority?**

Ensemble Learning

- **An ensemble model is a model that is a combination of several different models**
- **Usually, an ensemble is more accurate than all its constituent models**
- **Why?**
 - **Intuition:** "two know more than one"

Ensemble Learning. First approach: Voting

- Given n classifiers m_1, m_2, \dots, m_n
- Consider a new classifier M that, given a datum x , M computes $m_1(x), m_2(x), \dots$, counts the predictions and returns the most predicted class
- How well would M work?

Ensemble Learning. First approach: Voting

- To answer the question, let's make some assumptions:
 - All the classifiers m_1, \dots, m_n are equally accurate (with p being their accuracy)
 - The errors in the classification made by each classifier are independent
 - $P(m_j \text{ wrong} \mid m_k \text{ wrong}) = P(m_j \text{ wrong})$
- How well would M work?

Ensemble Learning. First approach: Voting

- **Example**

- **$p = 0.8$**

- **$n = 5$**

- **$P(\text{majority of 5 models is correct}) =$**

$$= \binom{5}{5} 0.8^5 0.2^0 + \binom{5}{4} 0.8^4 0.2^1 + \binom{5}{3} 0.8^3 0.2^2$$

$$= 0.942$$

- **Need empirical validation if assumptions not valid**

Ensemble Learning. First approach: Voting

- **Example**

- **$p = 0.8$**

- **$n = 5$**

- **$P(\text{majority of 5 models is correct}) =$**

$$= \binom{5}{5} 0.8^5 0.2^0 + \binom{5}{4} 0.8^4 0.2^1 + \binom{5}{3} 0.8^3 0.2^2$$

$$= 0.942$$

- **Need empirical validation if assumptions not valid**

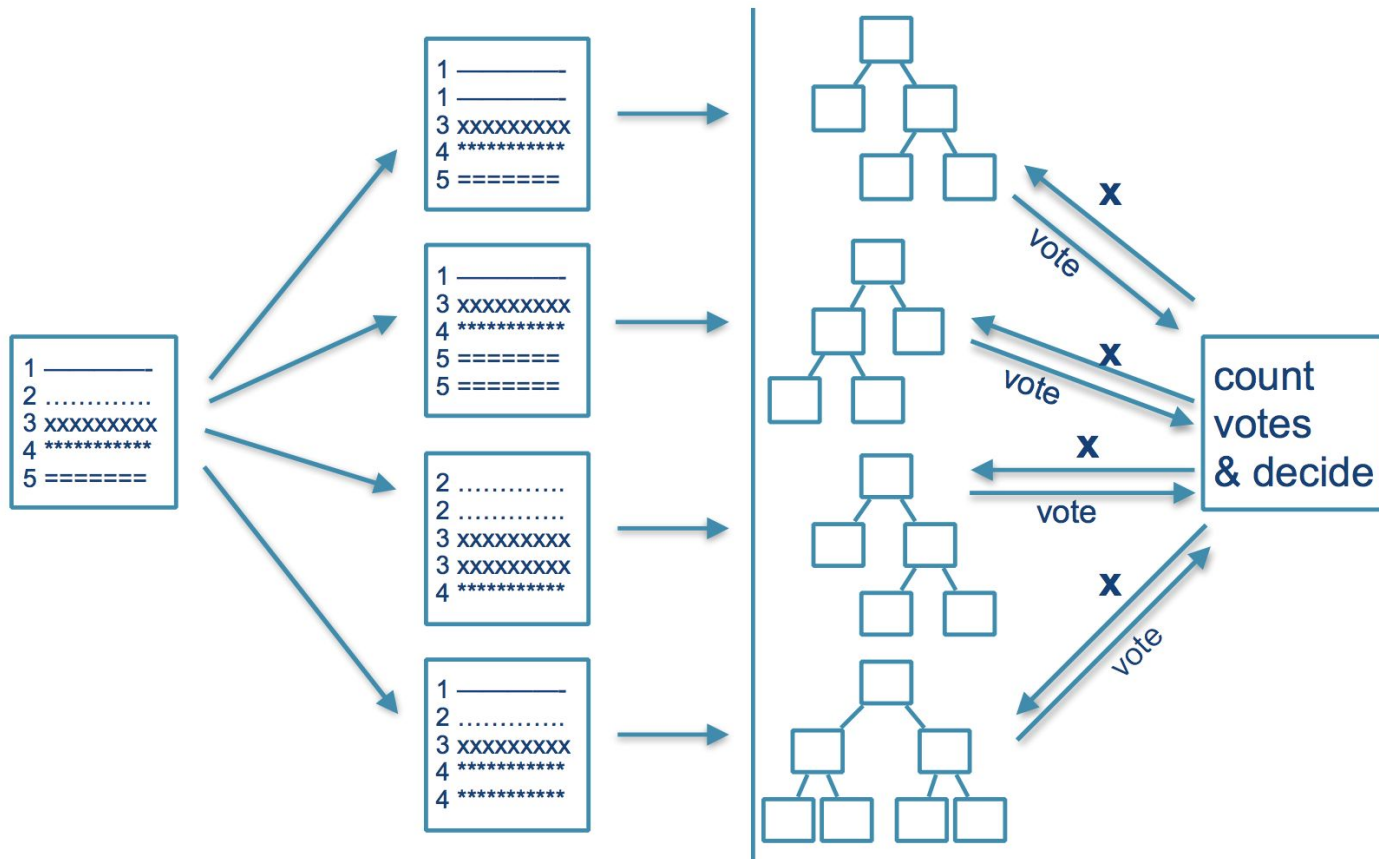
Ensemble Learning

- We know how to build decision trees
- How can we build different trees for the same data?
 - This may result on the same tree
- Do we introduce some variations?

Bagging

- Let D be the dataset
- Repeat k times:
 - Create D' from D by randomly selecting $|D|$ instances of D with replacement
 - Learn a new model m
- Return a model that selects the most frequent prediction among m_1, \dots, m_k predictions

Bagging



Bagging for Decision Trees

- **Bagging generally works well for unstable learners**
 - **A learner is unstable if small changes in the dataset can give very different resulting models**
 - **It turns out that decision tree learners are indeed unstable**
- **Disadvantage: learning k trees is k times as expensive as learning one tree**

Random Forests

- **Like bagging, with one improvement**
 - **For trees, ALL the features are considered to create a split node (inner node)**
 - **For random forests, at each node consider only M randomly chosen attributes (not all)**
 - **Usually take** $M = \sqrt{\text{number of attributes}}$

Random Forests

Common Steps

- **Build a random forest considering M attributes**
- **Predict the value of "Out-of-bag" samples using the random forest**
- **Estimate the accuracy**
- **Determine the optimal M (hyperparameter)**

Random Forests

- **Random Forests is one of the most efficient and most accurate learning methods to date (2008)** (Caruana+: An empirical evaluation of supervised learning in high dimensions. ICML 2008)
- **Easy to use with little parameter tuning**
- **Easy to debug, but, compared with Decision Trees, the model is less interpretable**

Demo Time (demo 6)

...

Boosting

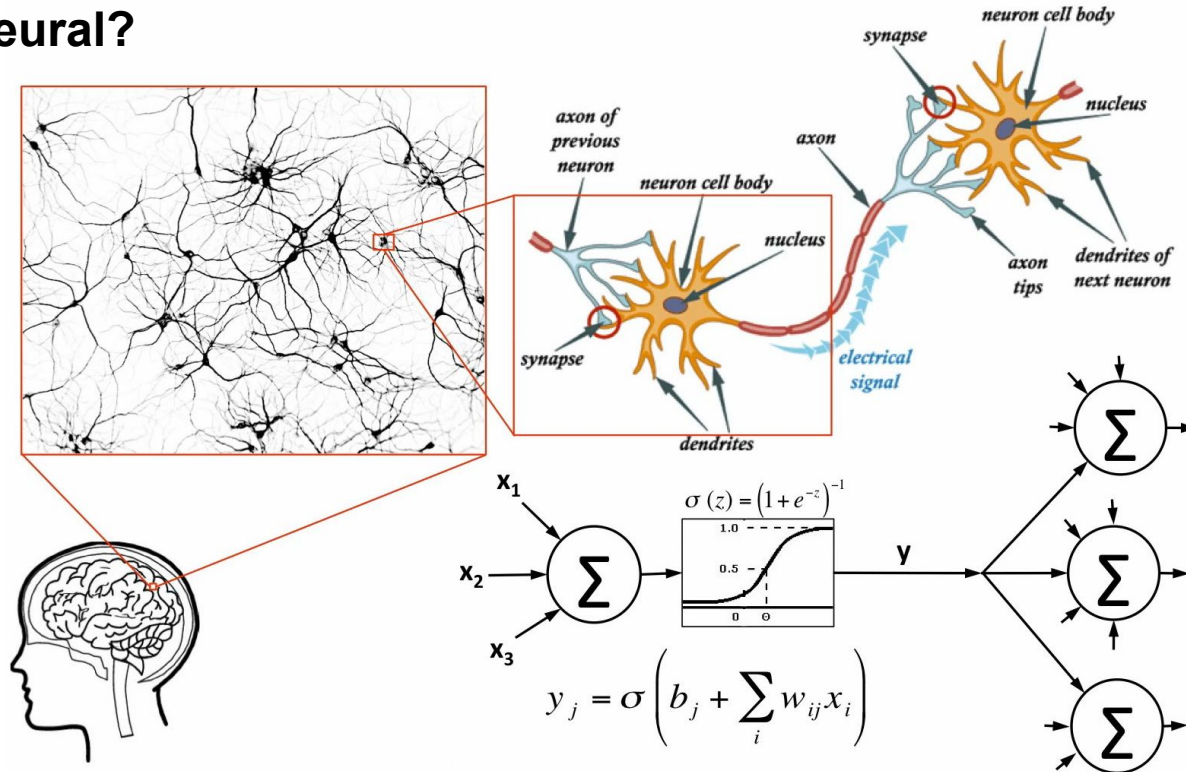
- **Bagging goal:** fit large trees to resampled versions of the training data, and classify by majority vote
- **Boosting:** fit large or small trees to "**reweighted**" versions of the training data and classify by **weighted** majority vote

Boosting

- **Each model, defines the features that the next model will focus on**
- **Uses bootstrapping like bagging, but here we weight each sample of data**
 - **Some samples will be used more frequently**
- **Process:**
 - Given a model, track the samples that are more "erroneous" and give them heavier weights (considered to be data that have more complexity and requires more steps)
 - Given a model, track the error rate so that better models are given more weights

Neural Networks Warm-up

Why neural?



Neural Networks Warm-up

Notation

Remember that we are given a set of labeled samples:

- **m training samples:** $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots (x^{(m)}, y^{(m)})\}$
- $X = [x^{(1)} x^{(2)} \dots x^{(m)}] \in \mathbb{R}^{n_x \times m}$
- $Y = [y^{(1)} y^{(2)} \dots y^{(m)}] \in \mathbb{R}^{1 \times m}$

Neural Networks Warm-up

Logistic Regression Review

Given x , we would like to find $\hat{y} = P(y = 1|x)$

What is the easiest way to transform x ?

Neural Networks Warm-up

Logistic Regression Review

Given x , we would like to find $\hat{y} = P(y = 1|x)$

What is the easiest way to transform x ?

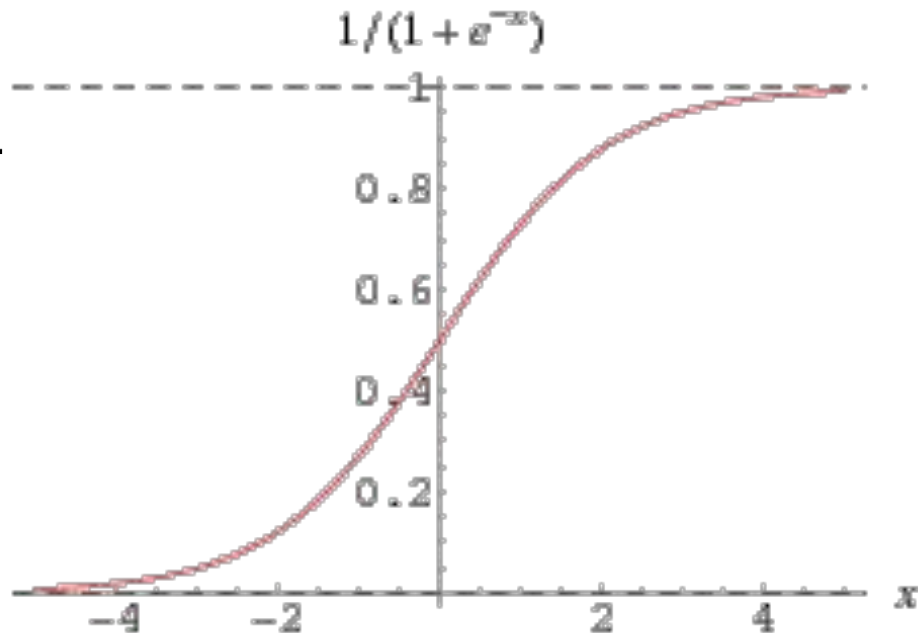
$$\hat{y} = w^T x + b$$

But we would like \hat{y} to be a probability: $0 \leq \hat{y} \leq 1$

Neural Networks Warm-up

Sigmoid function

$$f(x) = \frac{1}{1 + \exp(-x)}$$



Neural Networks Warm-up

Cost function (for Logistic Regression)

- **m training samples:** $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots (x^{(m)}, y^{(m)})\}$
- **want to find:** $\hat{y}^{(i)} \approx y^{(i)}$
- **Popular loss function:**

$$\mathcal{L}(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$$

Neural Networks Warm-up

Cost function (for Logistic Regression)

- **m training samples:** $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

- **want to find:** $\hat{y}^{(i)} \approx y^{(i)}$

- **Popular** loss function:

$$\mathcal{L}(\hat{y}, y) = -[y \log \hat{y} + (1 - y) \log (1 - \hat{y})]$$

- **Cost function:**

$$\mathcal{J}(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

Neural Networks Warm-up

Logistic Regression

$$z = w^T x + b$$

$$\hat{y} = \sigma(w^T x + b) \text{ (activation function; we usually use } a \text{ instead of } \hat{y})$$

$$\mathcal{L}(a, y) = -[y \log a + (1 - y) \log (1 - a)]$$

The goal is to minimize \mathcal{J} , to do so we compute

$$\frac{\partial}{\partial w_1} \mathcal{J}(w, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial w_1} \mathcal{L}(a^{(i)}, y^{(i)})$$

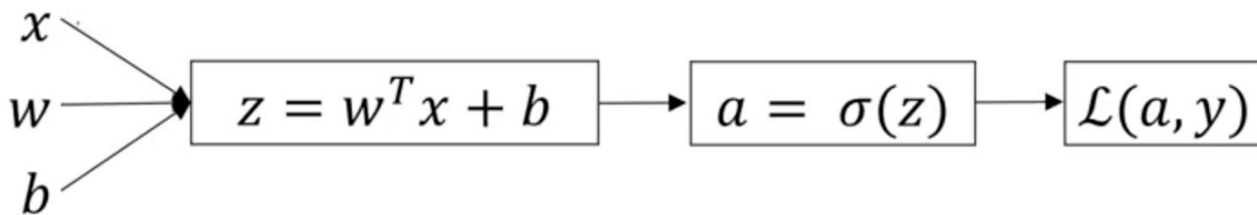
Whiteboard Time

...

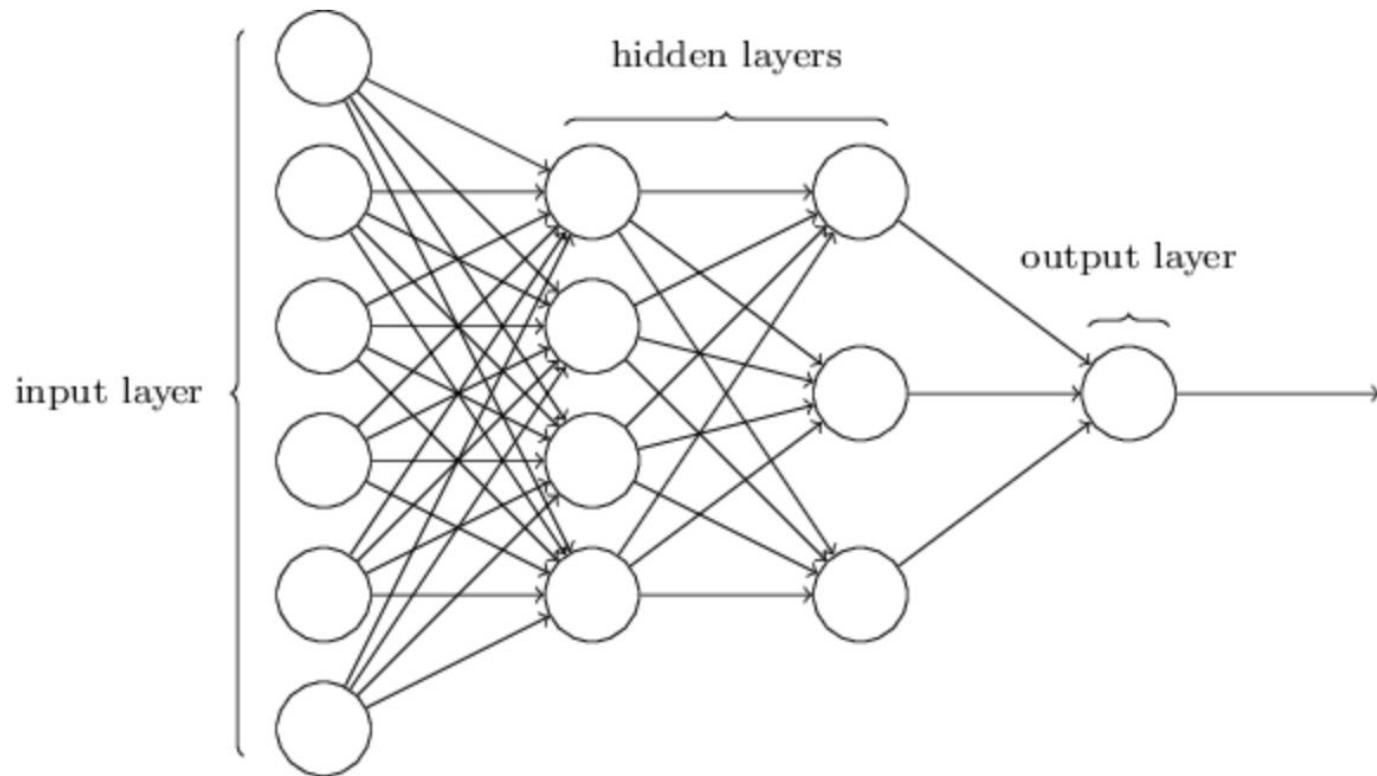
Computation graph for backpropagation
&
Logistic Regression minimization

Neural Networks Warm-up

Logistic Regression Gradient Descent



Neural Networks



Neural Networks

Activation Functions

- Sigmoid function (as in Logistic regression)
- **tanh:** $\frac{\exp(x) - \exp(-x)}{\exp(x) + \exp(-x)}$
- **Rectified Linear Unit (ReLU):** $\max(0, x)$
- Leaky ReLU

Demo Time (demo 7)

...

Demo Time (demo 7)



The screenshot shows the scikit-learn documentation page for "1.17. Neural network models (supervised)". The page has a navigation bar with links for Home, Installation, Documentation, and Examples. A search bar is also present. On the left sidebar, there are links for Previous, Next, and Up versions, as well as a link to cite the software. The main content area features a warning message highlighted with a red box, stating that the implementation is not intended for large-scale applications and lacks GPU support. Below the warning, the section "1.17.1. Multi-layer Perceptron" is visible.

scikit-learn

Home Installation Documentation Examples

Google Custom Search

Previous 1.16. Probabl... Next 2. Unsupervis... Up 1. Supervised...

scikit-learn v0.19.1
[Other versions](#)

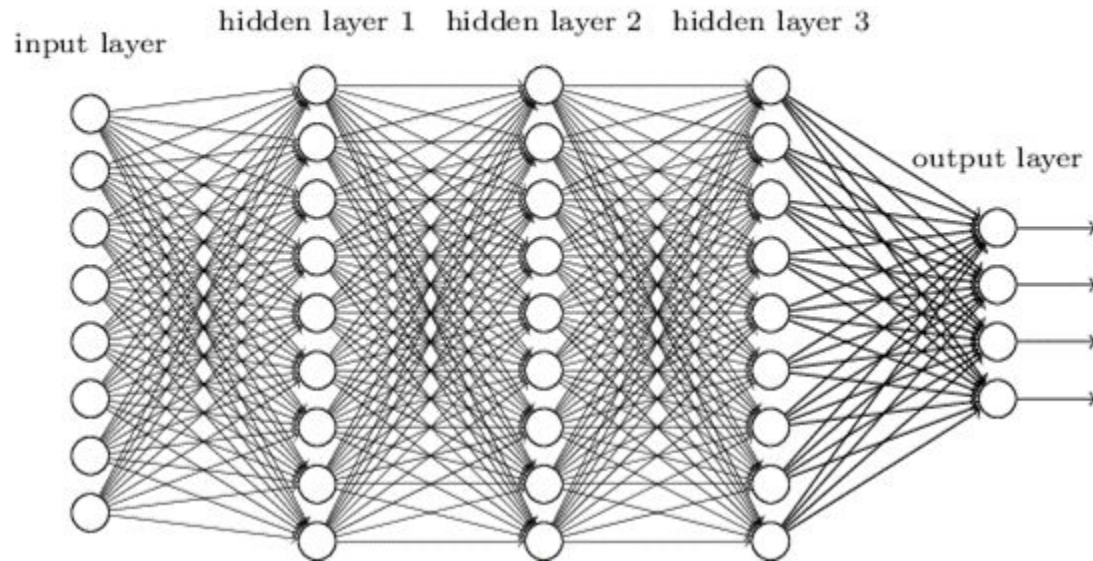
Please [cite us](#) if you use the software.

1.17. Neural network models (supervised)

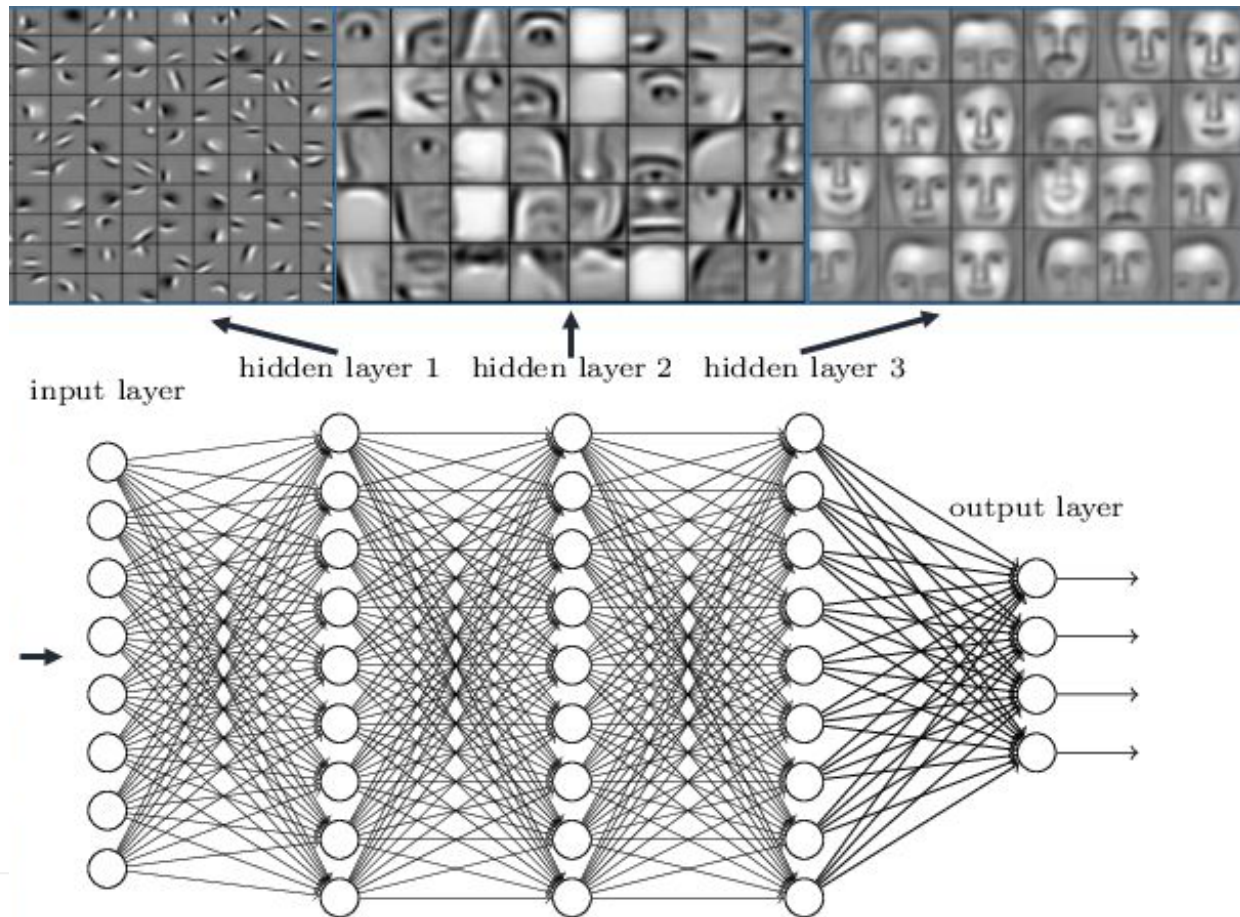
Warning: This implementation is not intended for large-scale applications. In particular, scikit-learn offers no GPU support. For much faster, GPU-based implementations, as well as frameworks offering much more flexibility to build deep learning architectures, see [Related Projects](#).

1.17.1. Multi-layer Perceptron

Deep Neural Networks



Deep Neural Networks



Deep Neural Networks

How to split the data?

- **Train / Test / Validation**



Deep Neural Networks

How to split the data?

- **Train / Test / Validation**



- **Now?**

Deep Neural Networks

How to split the data?

- **Train / Test / Validation**



- **Now?**
 - **Too much data (> 10.000.000 samples)**



Deep Neural Networks

How to split the data?

- Train / Test / Validation



- Now?
 - Too much data (> 10.000.000 samples)



**Make sure your train/ test /
validation come from the
same distribution**

Deep Neural Networks

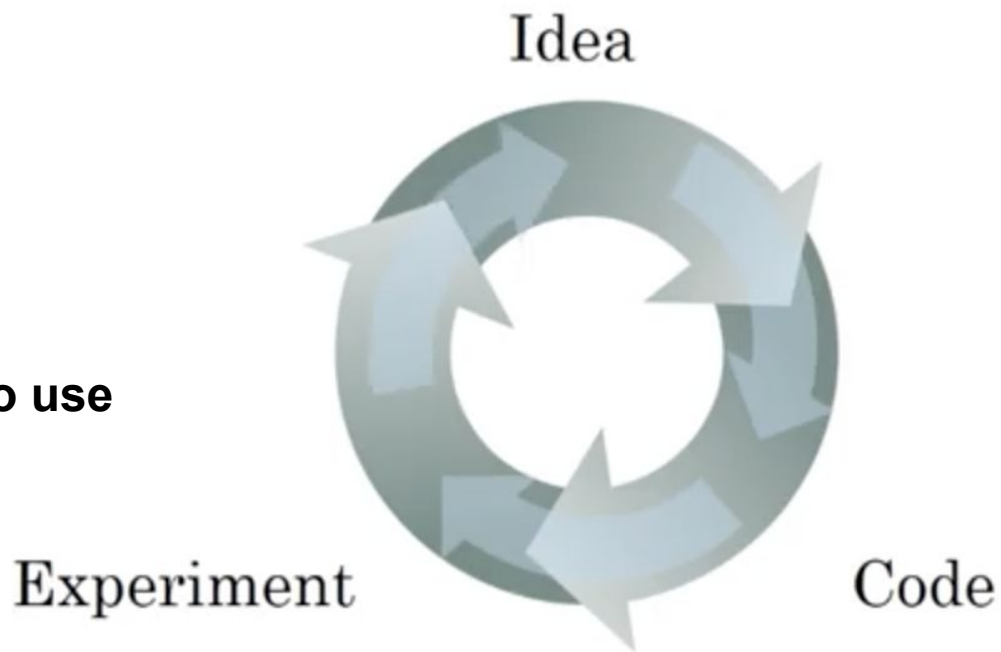
Rule of thumb to deal with bias/variance?

- **High Bias:**
 - **Bigger Network**
 - **Different Network Architecture**
- **High Variance:**
 - **More data**
 - **Regularization**
 - **Different Network Architecture**

Neural Networks

How to decide the size?

- # hidden layers
- # hidden units
- What activation function to use



Neural Networks Regularization

Logistic regression:

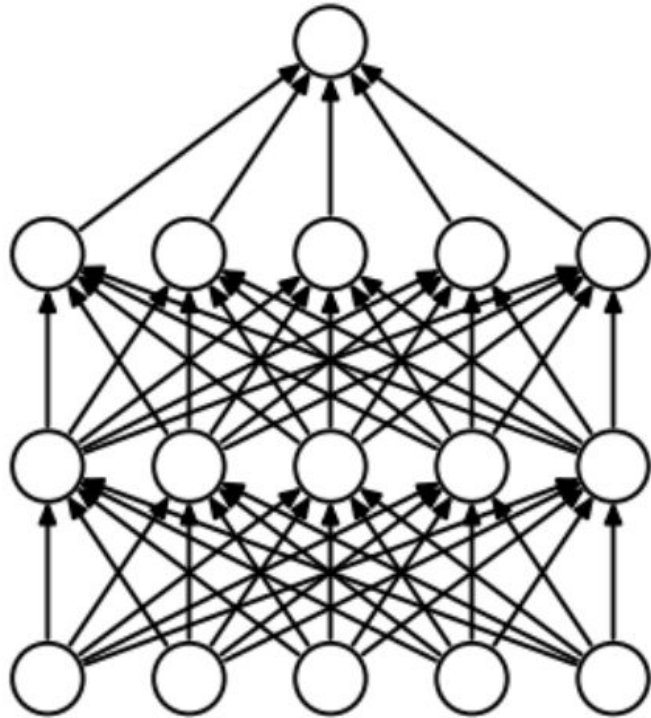
$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

Neural network:

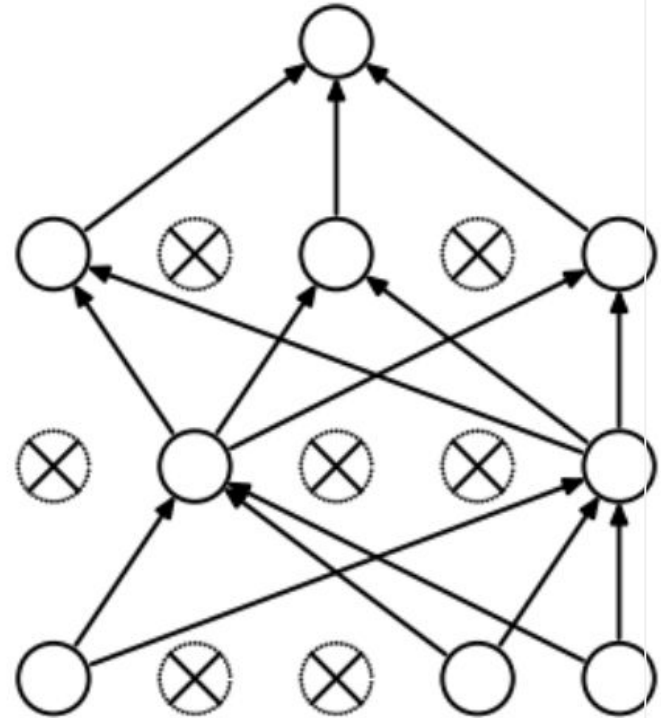
$$h_{\Theta}(x) \in \mathbb{R}^K \quad (h_{\Theta}(x))_i = i^{th} \text{ output}$$

$$J(\Theta) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_{\Theta}(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - (h_{\Theta}(x^{(i)}))_k) \right] \\ + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_{ji}^{(l)})^2$$

Neural Networks Dropout



(a) Standard Neural Net



(b) After applying dropout.

Neural Networks Momentum

- **Mini-batch Gradient Descent**
- **Gradient Descent optimization**
 - **with Momentum**
 - **RMSprop**
 - **Adam**

First Project