

Aprendizaje Supervisado

2da Parte

...

Franco Luque - Matías Marenchino

Agenda

- Repaso
Aprendizaje Supervisado / SVMs / Ensembles / NNs
- Resolución Laboratorio 1
- Estrategias para Machine Learning
- Caso de Estudio: Clasificación de Texto

Repaso

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

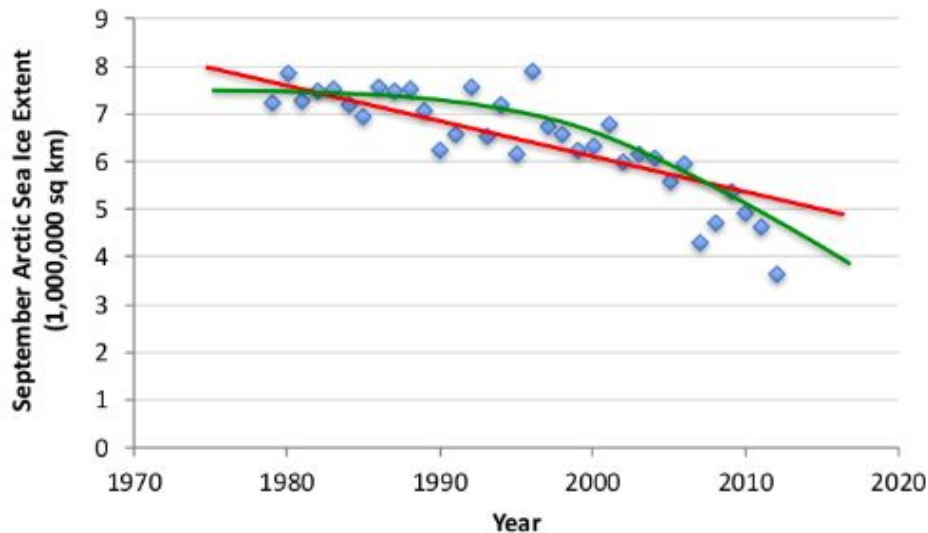
- **Problem:** classify an email as SPAM or not

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

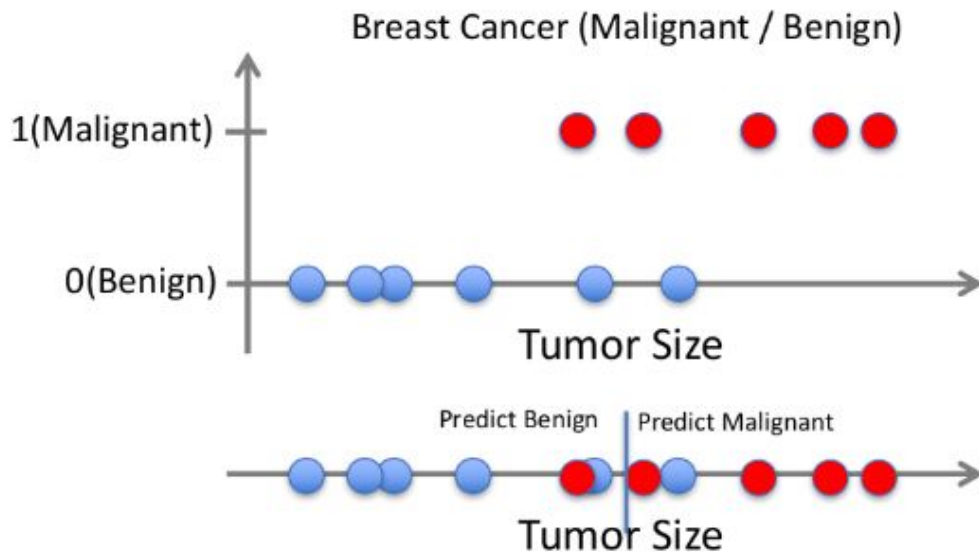
- Datos $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$
- Aprender una $f(x)$ que permita predecir y a partir de x
 - Si y está en $\mathbb{R}^n \rightarrow$ **regresión**



Slides from the previous course

Aprendizaje supervisado: clasificación

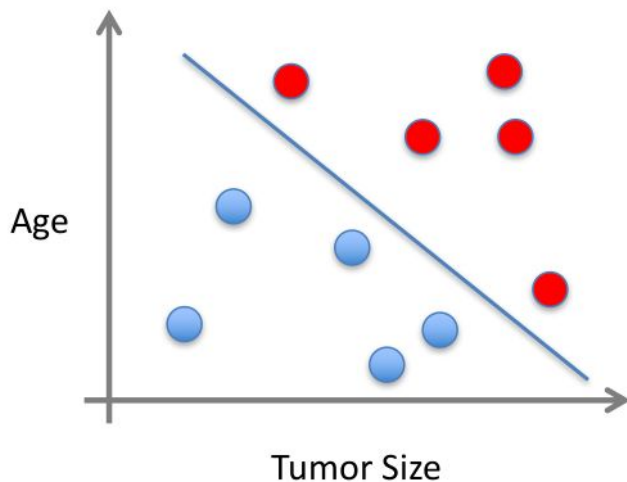
- Datos $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$
- Aprender una $f(x)$ que permita predecir y a partir de x
 - Si y es categórica \rightarrow **clasificación**



Slides from the previous
course

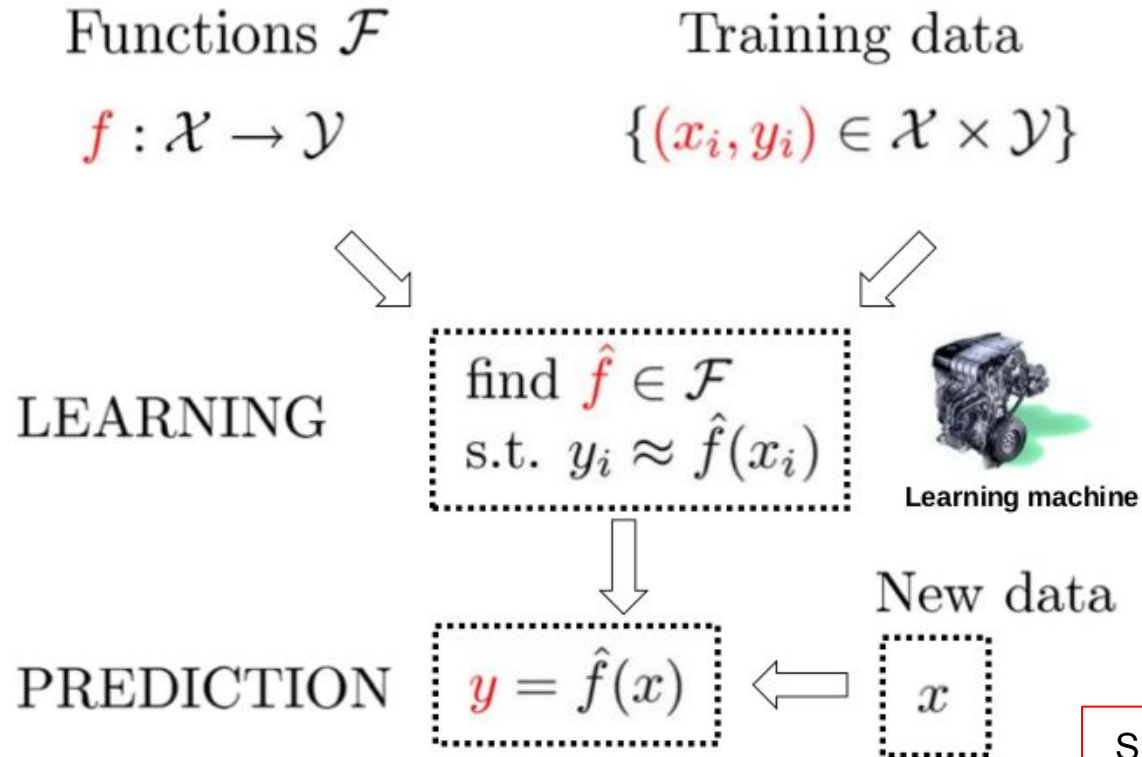
Supervised Learning

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



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

Aprendizaje supervisado

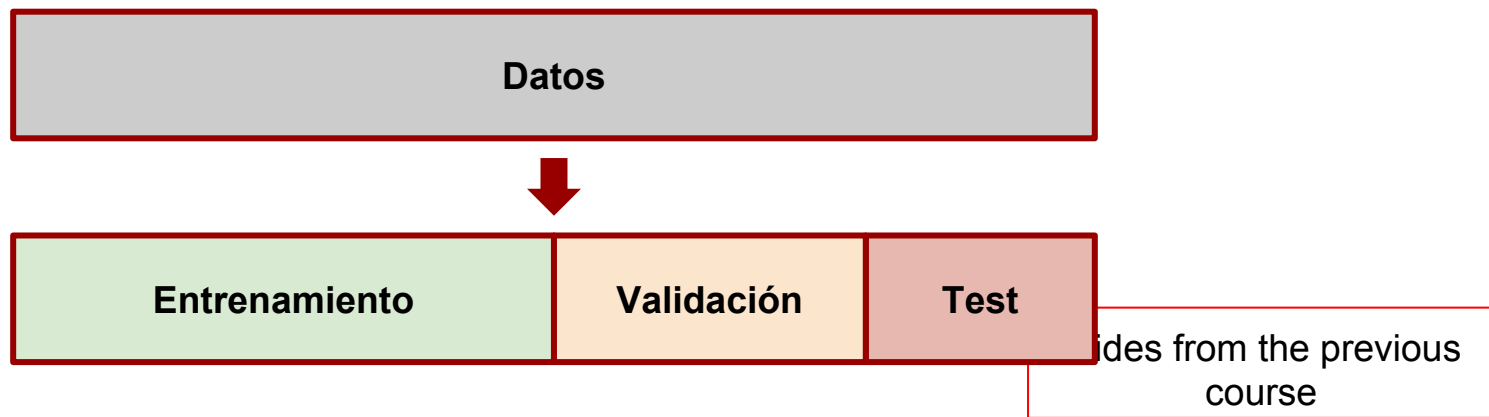


Slides from the previous
course

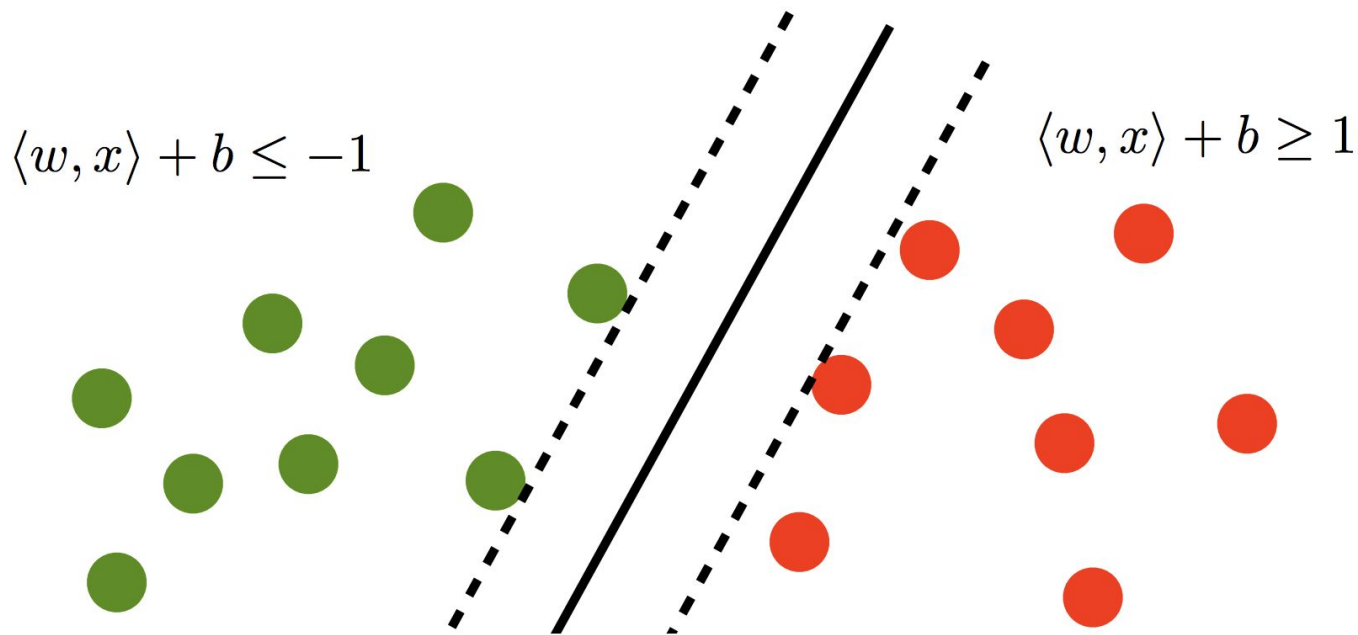
Elección de *hiperparámetros*

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)



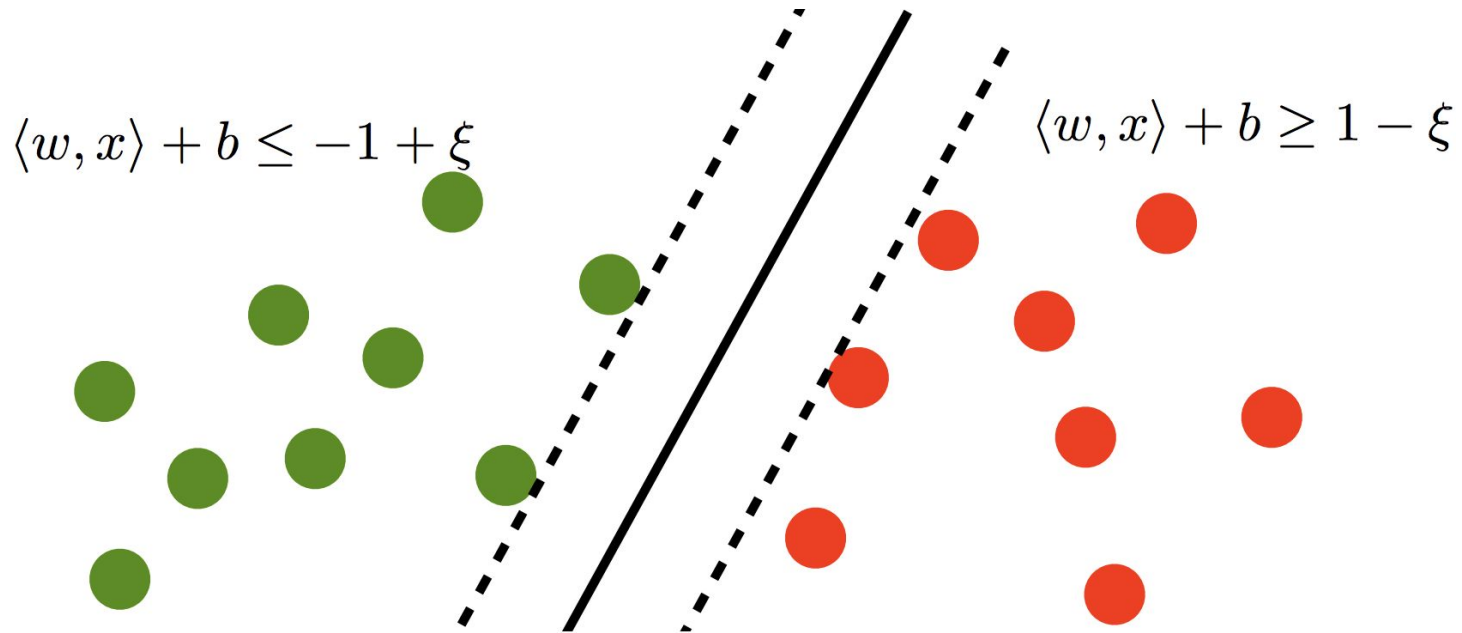
Support Vector Machines



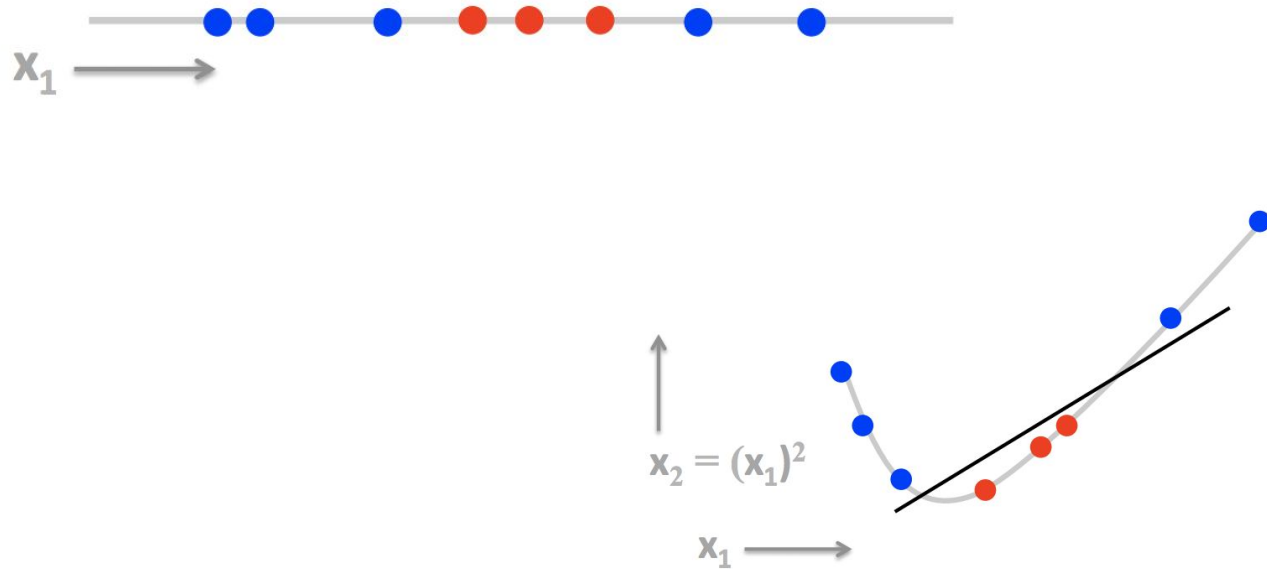
linear function

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

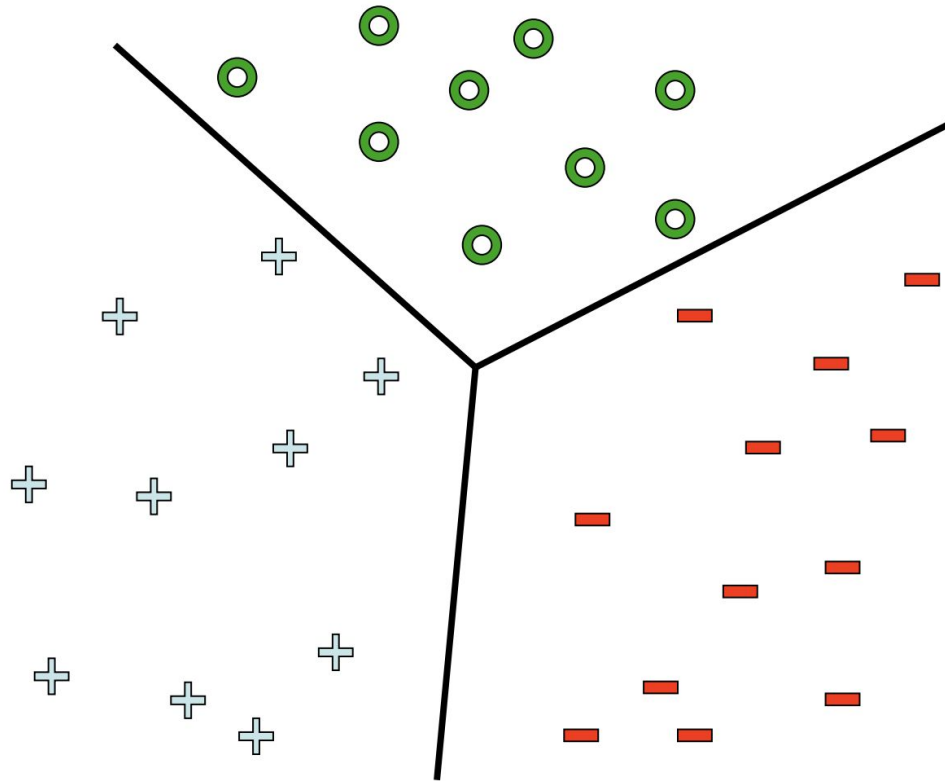
SVMs: slack variables



SVMs: Kernels



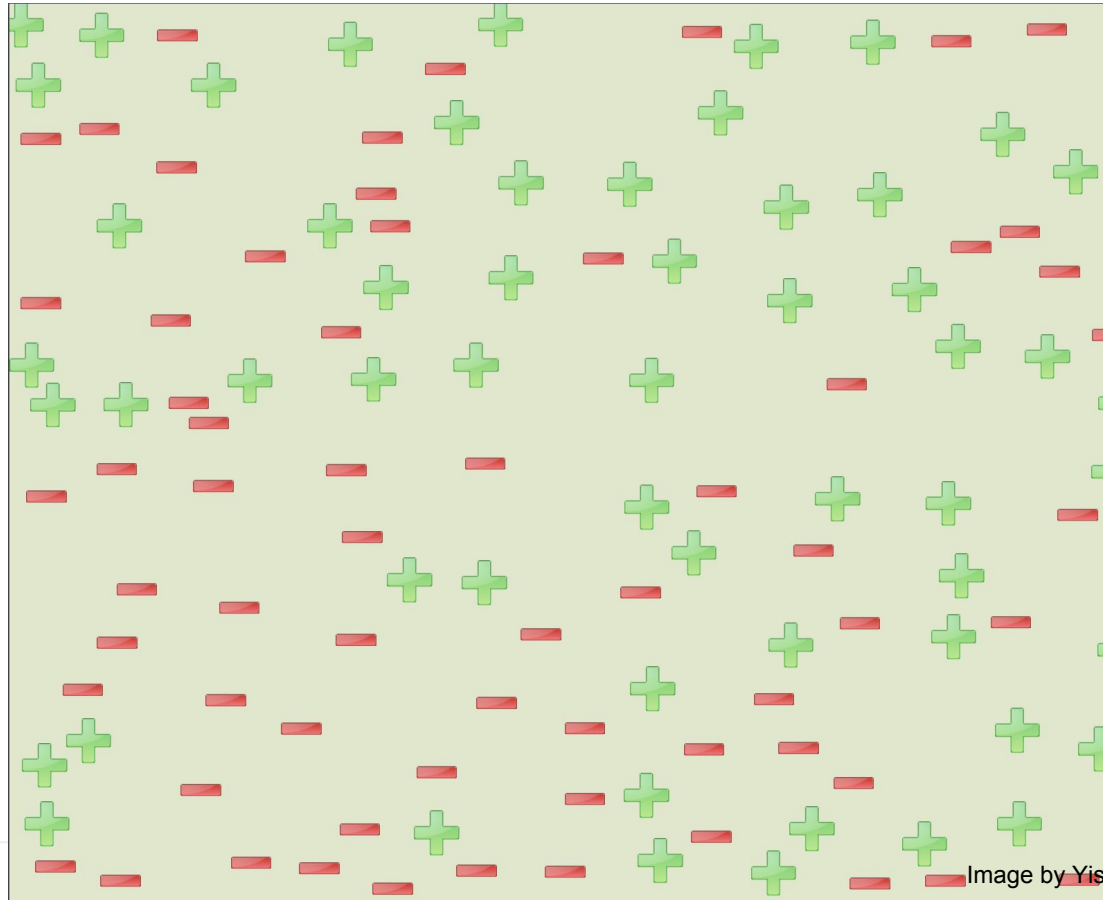
Multiclass SVMs



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.

Decision Trees (review)



Decision Trees (review)



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"

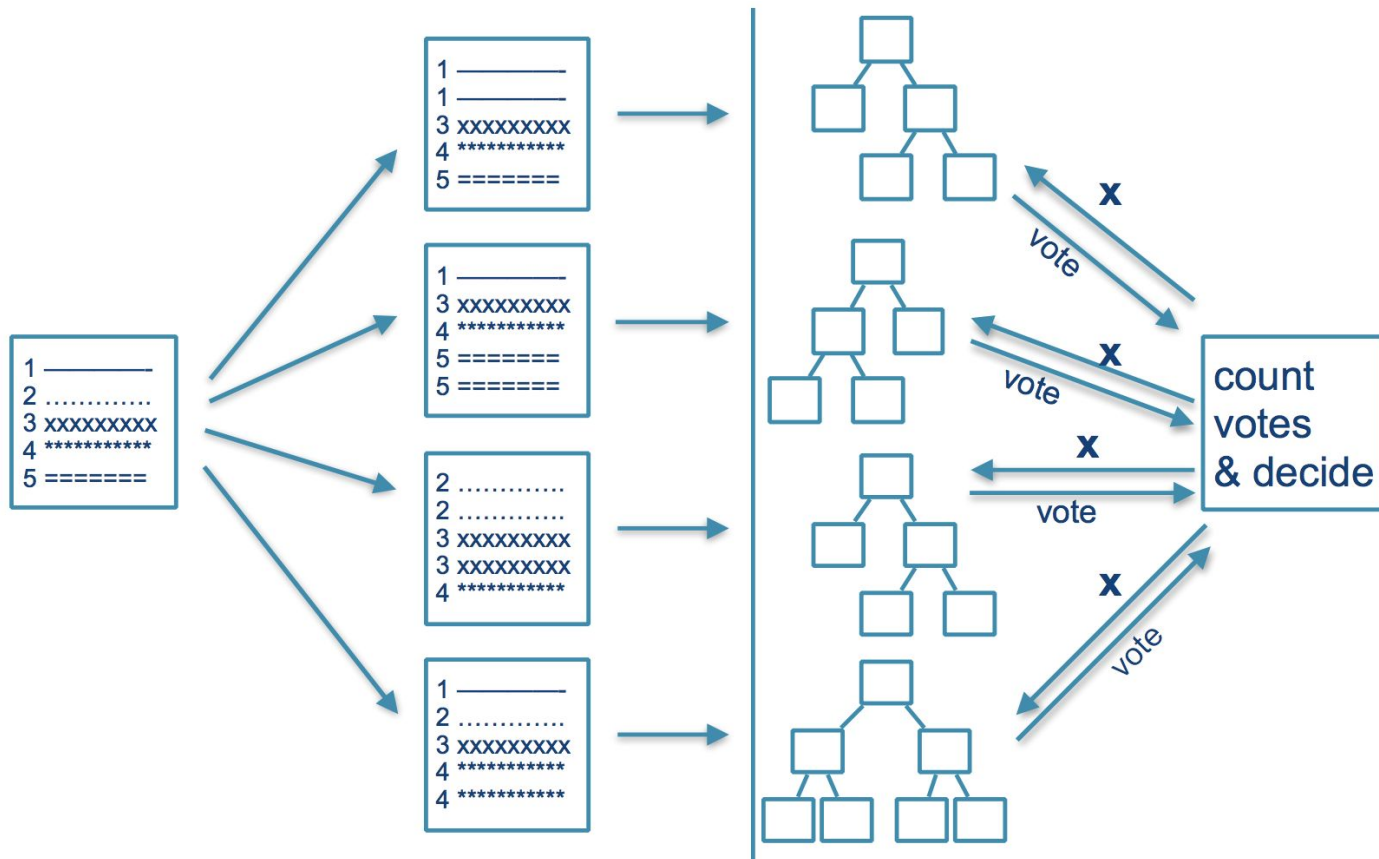
Ensemble Learning. First approach: Voting

- 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?

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**

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

Logistic Regression Review

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

What is the easiest way to transform x ?

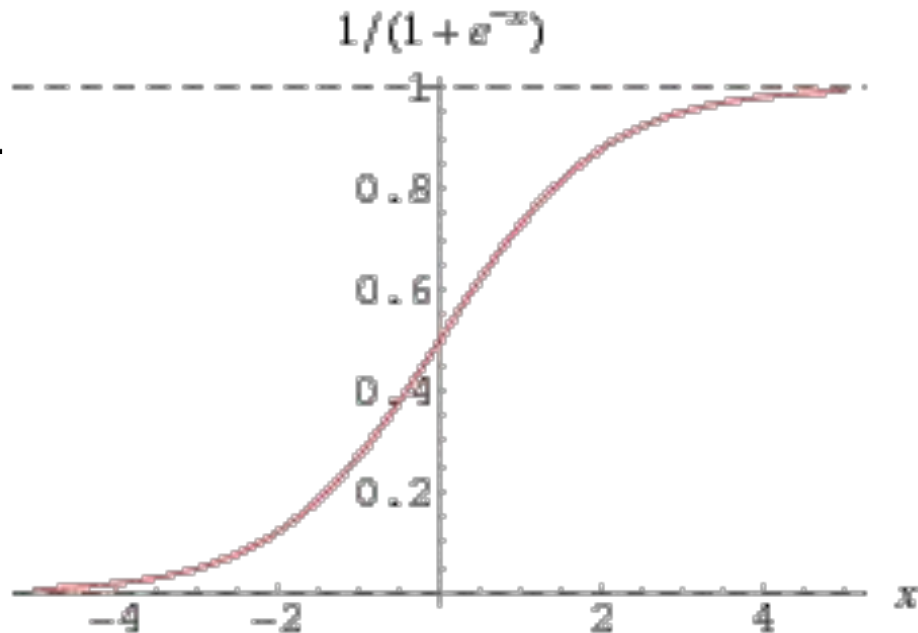
$$\hat{y} = w^T x + b$$

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

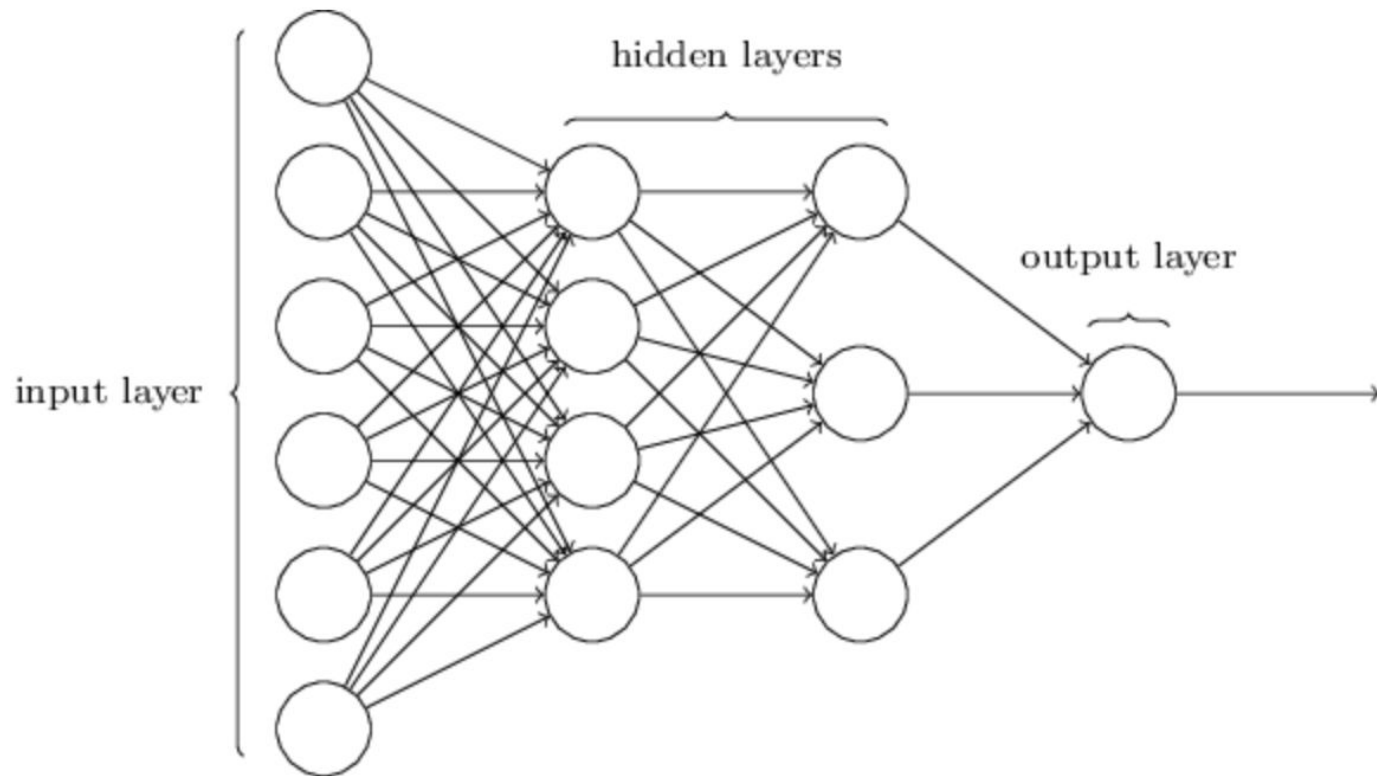
Neural Networks Warm-up

Sigmoid function

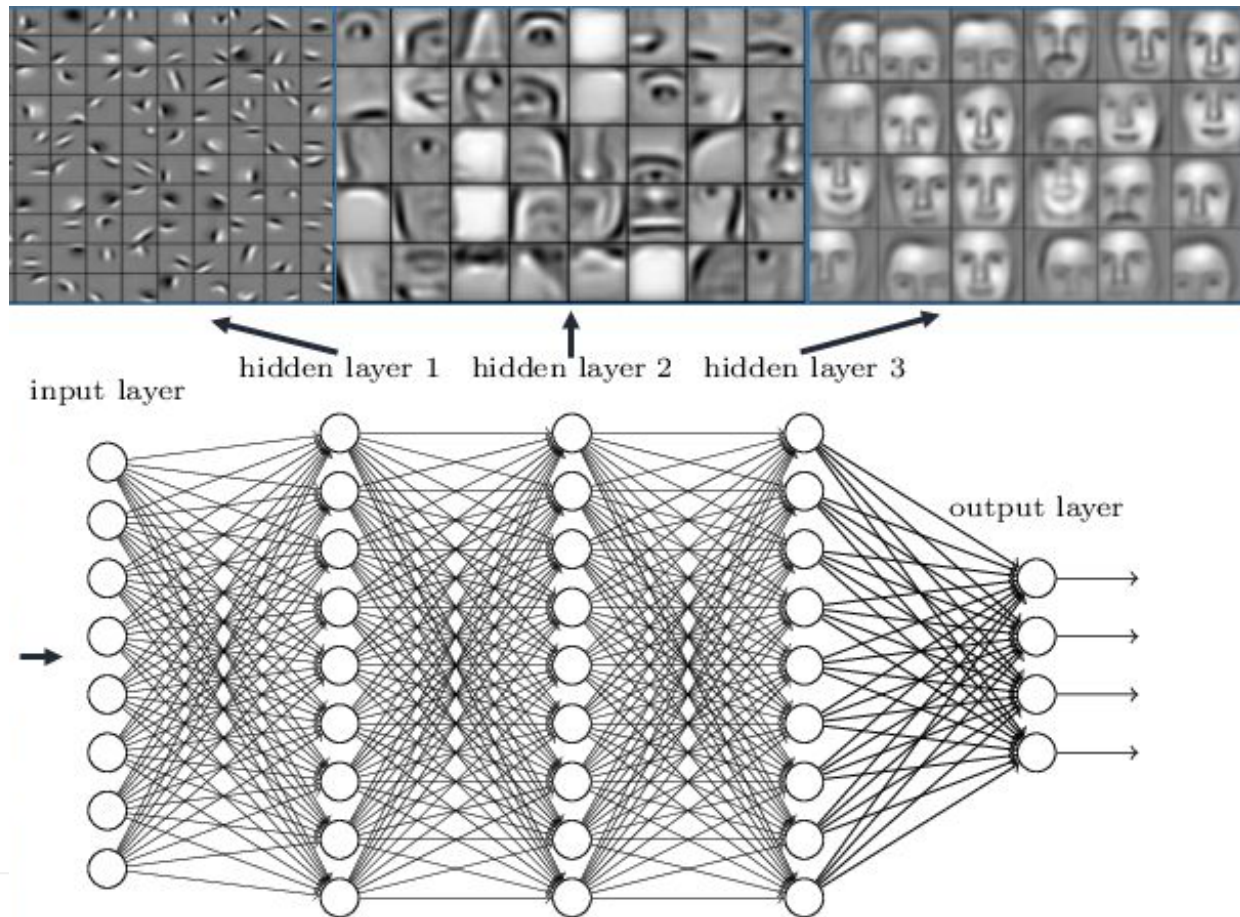
$$f(x) = \frac{1}{1 + \exp(-x)}$$



Neural Networks



Deep Neural Networks



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

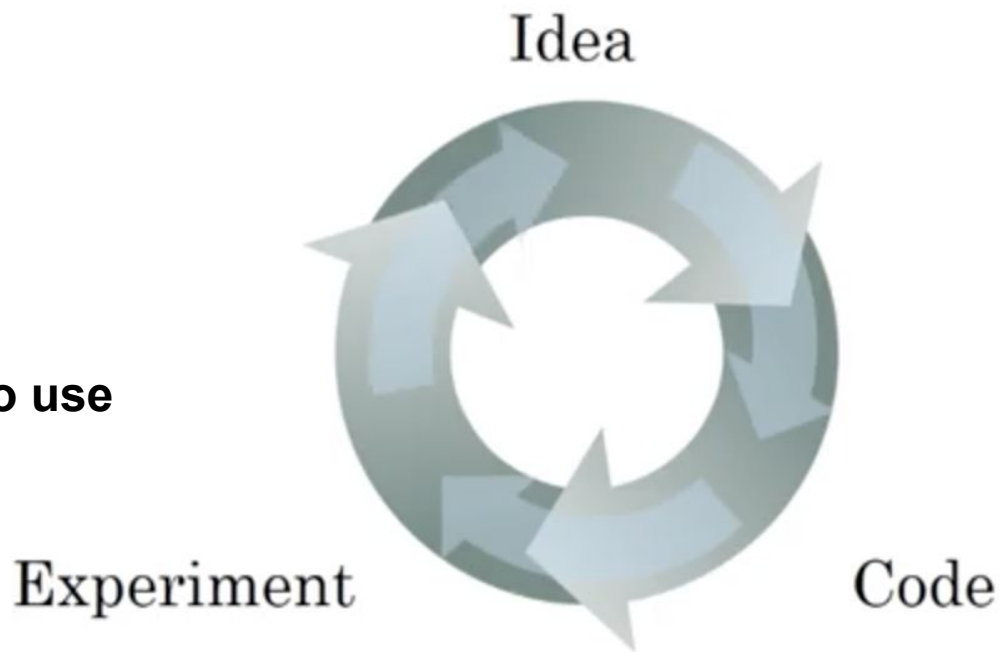
Rule of thumb to deal with bias/variance?

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

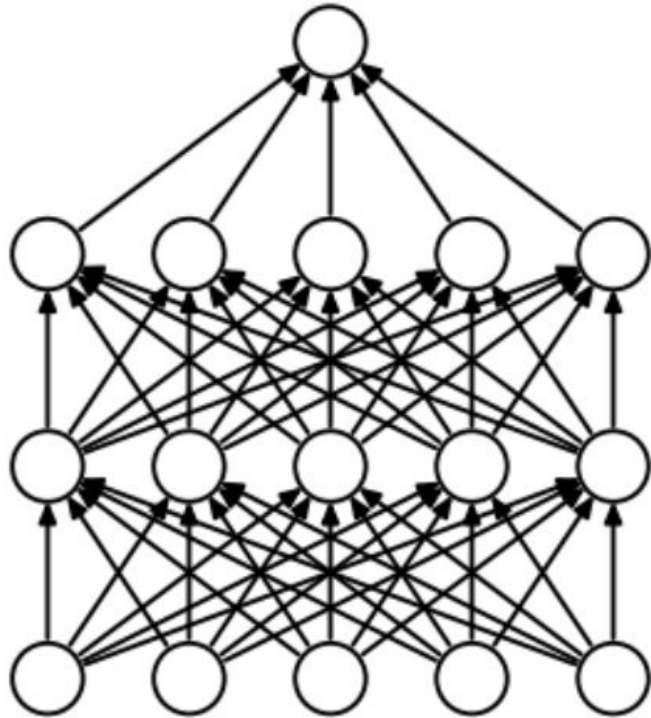
Neural Networks

How to decide the size?

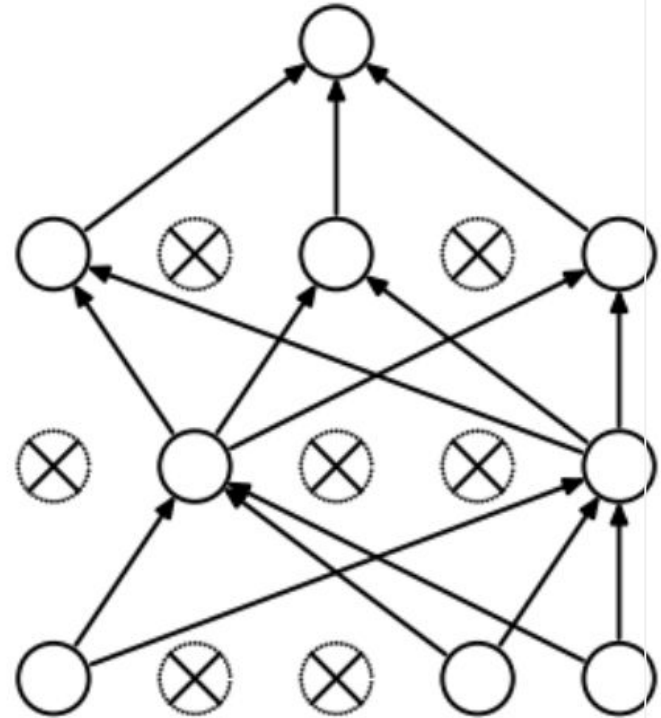
- # hidden layers
- # hidden units
- What activation function to use



Neural Networks Dropout



(a) Standard Neural Net



(b) After applying dropout.