

Numerical Computation and Machine Learning Basics

<https://github.com/timestocome/DeepLearning-Talks>

Numerical computing resources:

Numerical Recipes in C/C++, Fortran 77/90

<http://numerical.recipes/oldverswitcher.html>

Accelerated Numerical Analysis Tools with GPUs

<https://developer.nvidia.com/how-to-numerical-analysis>

Overflow and Underflow

Overflow -+/- infinity, NaN, wrapping

Underflow - numbers fade to zero

vanishing gradients [clip gradients, add epsilon]

exploding gradients [clip gradients, use regularization]

- Pentium FDIV Bug
- GPU truncation causing non-linearity
- <https://blog.openai.com/nonlinear-computation-in-linear-networks/>

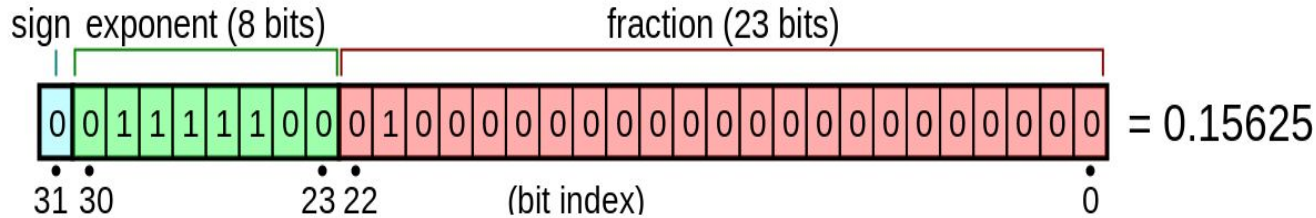
Most ML libraries will handle NaN, -/+ inf, div by zero for you

Wrapping: use a 16 bit int to count in a loop and it'll wrap at 65,536 - 1

- don't use loops
- use a very large loop index

IEEE Standard

$(-1 * s) \times (\text{fraction}) \times (2 ^ \text{exponent})$



(* fraction is usually named mantissa)

Size of ints, longs, floats, doubles on your computer

<https://www.programiz.com/c-programming/examples/sizeof-operator-example>

http://cstl-csm.semo.edu/xzhang/Class%20Folder/CS280/Workbook_HTML/FLOATING_tut.htm

Poor Conditioning

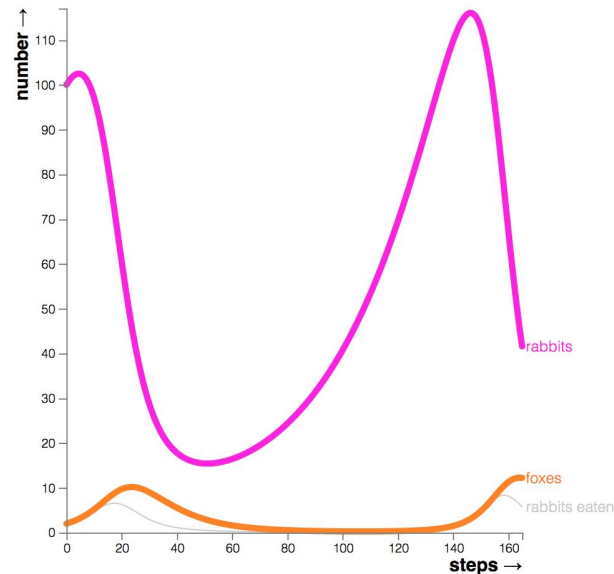
Conditioning is how quickly a function changes

Even a small learning rate increases the error causing the network to become stuck

In a small network this is a local minima, in a deep network it's a saddle point

Fix: use an adaptive learning rate (Adagrad, AdaDelta, RMSProp, Adam, Momentum ..

Adaptive learning rates are built into Theano, TensorFlow ...

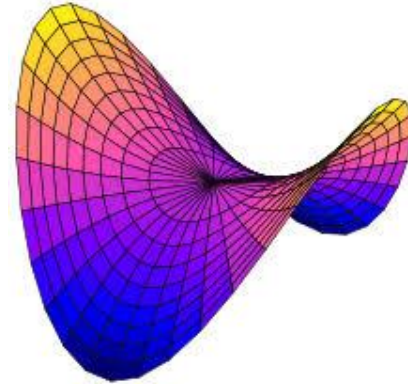


Gradient Based Optimization

Goal is to minimize or maximize a cost function (loss function, error function)

Using a small learning rate, take a small step in the direction of steepest decrease

Gradient = 0 at minimum, max, saddle point



Derivative

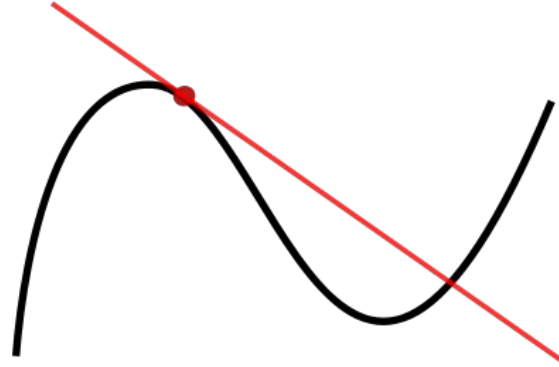
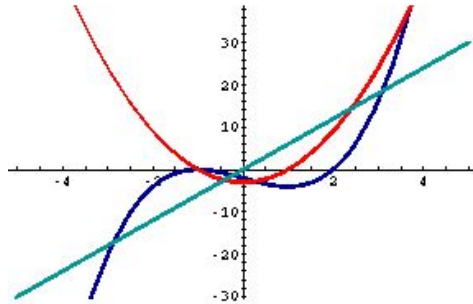
Measures the tangent line on a slope which gives you the direction of greatest increase. In neural networks we use the negative of it to get steepest descent.

The 2nd derivative can tell you if you are at a max or min

$f'' < 0$ max

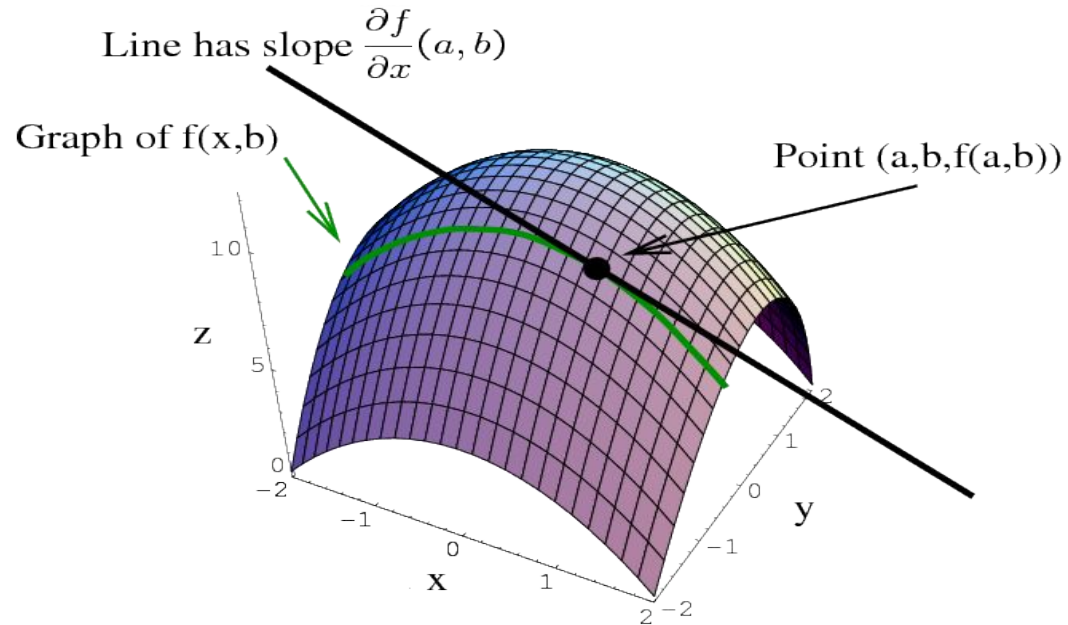
$f'' = 0$ no information

$f'' > 0$ minimum



Partial Derivative

Derivatives calculated on vectors (gradient)



Jacobian Matrix

Take the first derivative of a vector and output a vector to find steepest descent

$$\frac{\partial (x, y)}{\partial (u, v)} = \begin{vmatrix} \frac{\partial x}{\partial u} & \frac{\partial x}{\partial v} \\ \frac{\partial y}{\partial u} & \frac{\partial y}{\partial v} \end{vmatrix} = \frac{\partial x}{\partial u} \frac{\partial y}{\partial v} - \frac{\partial x}{\partial v} \frac{\partial y}{\partial u}$$

Hessian Matrix

Used to find the 2nd derivative of a vector

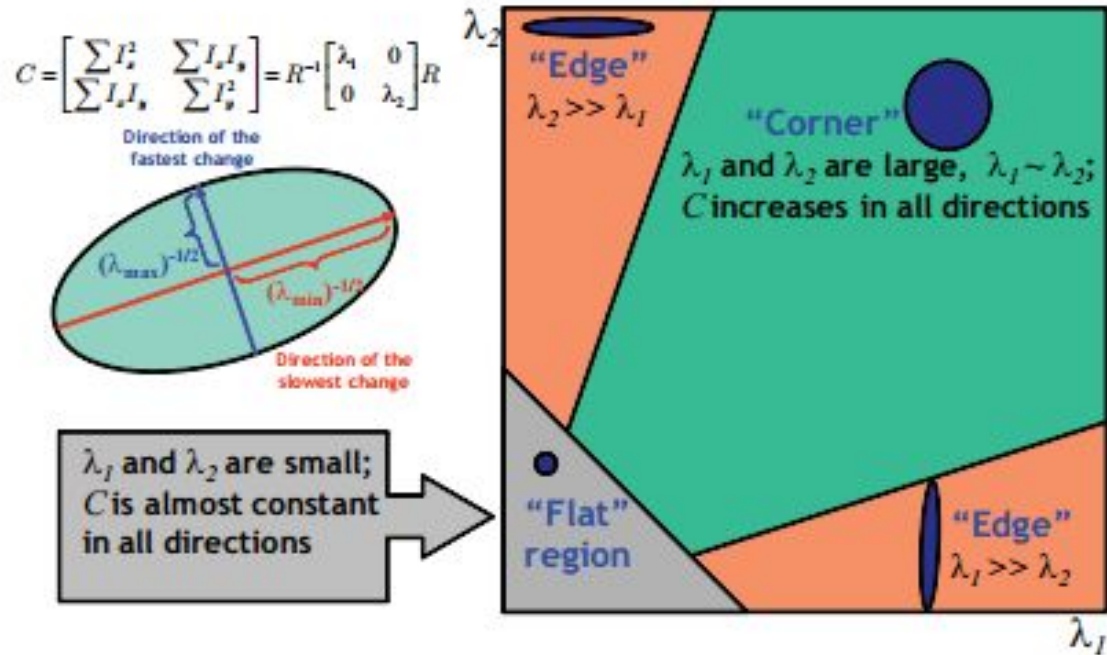
negative -> maxes negative curvature is down,

positive -> mins curvature is up

zero -> flat space or a saddle point

$$\text{Gradient} = \begin{bmatrix} \frac{\delta H}{\delta z_1} \\ \frac{\delta H}{\delta z_2} \\ \frac{\delta H}{\delta z_3} \end{bmatrix}, \text{Hessian} = \begin{bmatrix} \frac{\partial^2 H}{\partial z_1^2} & \frac{\partial^2 H}{\partial z_1 \partial z_2} & \frac{\partial^2 H}{\partial z_1 \partial z_3} \\ \frac{\partial^2 H}{\partial z_2 \partial z_1} & \frac{\partial^2 H}{\partial z_2^2} & \frac{\partial^2 H}{\partial z_2 \partial z_3} \\ \frac{\partial^2 H}{\partial z_3 \partial z_1} & \frac{\partial^2 H}{\partial z_3 \partial z_2} & \frac{\partial^2 H}{\partial z_3^2} \end{bmatrix}$$

Eigenvalues of Hessian



Newton's Method

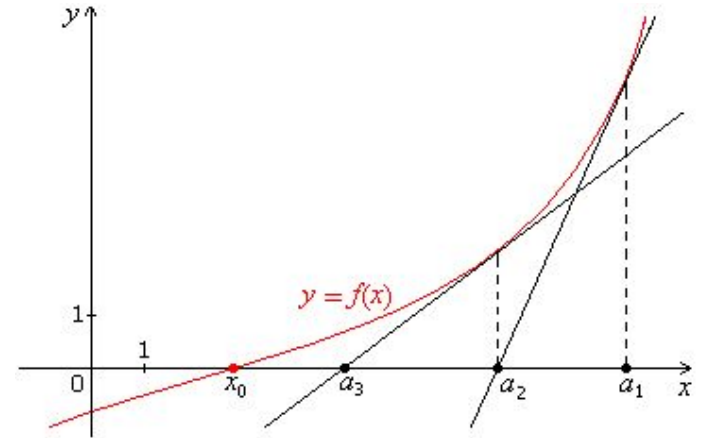
Approximate $f(x) = 0$

Keep guessing, split the difference between guesses until $f(x) = 0$

Sometimes fails spectacularly, extremely dependent on step size

Code example:

<http://www.aip.de/groups/soe/local/numres/bookcpdf/c9-4.pdf>



Constrained Optimization

Means there is more than one cost in the problem

i.e. Find production number that maximizes income using only factory available

Or maximize $f(x)$ subject to $g(x)$

Often just add $g(x)$ to cost function as a penalty, but a large penalty can create ill conditioned Hessians causing the network to fail to converge

Lagrangian Multiplier

solve: $\max \log(x) + \log(y)$ given: $x + y = m$

Combine equations, set $g(x)$ to zero $\max \log(x) + \log(y) + L(m - x - y)$

Set derivatives = 0 $1/x - z = 0, 1/y - z = 0, m - x - y = 0$

Remove L from x, y $1/x - z = 1/y - z \implies x = y$

2 equations, 2 unknowns $x = y = m/2$

Solve for Lagrangian multiplier $z = m/2$

<http://home.uchicago.edu/~vlima/courses/econ201/pricetext/chapter2.pdf>

http://adl.stanford.edu/aa222/Lecture_Notes_files/constrainedOptimization.pdf

Karush-Kuhn-Tucker Approach

- Combining Lagrange conditions for equality and inequality constraints yields *KKT conditions* for general

problem: $\min_{\mathbf{x}} f(\mathbf{x})$

s. t. $\mathbf{g}(\mathbf{x}) \leq \mathbf{0}$

$\mathbf{h}(\mathbf{x}) = \mathbf{0}$

Lagrangian:

$$L = f(\mathbf{x}) + \boldsymbol{\mu}^T \mathbf{g}(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{h}(\mathbf{x})$$

\Rightarrow

$$\frac{\partial L}{\partial \mathbf{x}} = \frac{\partial f}{\partial \mathbf{x}} + \sum \mu_i \frac{\partial g_i}{\partial \mathbf{x}} + \sum \lambda_i \frac{\partial h_i}{\partial \mathbf{x}} = \mathbf{0} \quad (\text{optimality})$$

and

$$\mathbf{g} \leq \mathbf{0}, \quad \mathbf{h} = \mathbf{0} \quad (\text{feasibility})$$

$$\boldsymbol{\lambda} \neq \mathbf{0}, \quad \boldsymbol{\mu} \geq \mathbf{0}, \quad \mu_i g_i = 0 \quad (\text{complementarity})$$

Linear Least Squares

Solve: $f(x) = \frac{1}{2} \|Ax - b\|^2$

Gradient: $df(x) = A.T(Ax - b)$

while $\|A.TAx - A.Tb\| > \text{acceptable_error}$:

$x = x - \text{learning_rate}(A.TAx - A.Tb)$

* do not use Newton's Method

Linear Least Squares with Constraint

Solve: $f(x) = \frac{1}{2} \|Ax - b\|^2$

Constraint: $g(x) = x.Tx \leq 1$

Combine the equations: $\frac{1}{2}\|Ax-b\|^2 + z(x.Tx - 1)$

set derivative = 0 $A.TAx - A.Tb + 2zx = 0$

Solve for Lagrangian Multiplier $z = (A.TAx - A.Tb) * x/2$

and the constraint $dL/dz = x.Tx - 1$

Resources

Bay Area Deep Learning School 2016

<https://www.youtube.com/watch?v=eyovmAtUx0> and <https://www.youtube.com/watch?v=9dXiAecyJrY>

Deep Learning Summer School Montreal

<https://www.youtube.com/playlist?list=PL5bqlc6XopCbb-FvnHmD1neVIQKwGzQyR>

http://videolectures.net/deeplearning2017_montreal/

Neural Networks for Machine Learning (Hinton Lectures 78 videos)

https://www.youtube.com/playlist?list=PLoRI3Ht4JOcdU872GhiYWf6jwrk_SNhz9

Must Know Tricks in Deep Neural Networks

<http://lamda.nju.edu.cn/weixs/project/CNNTricks/CNNTricks.html>

Bayesian vs Frequentist

Bayesian

Prior beliefs

Fixed parameters

$$\mu = mx + b$$

$$y = N(\mu, \text{std})$$

N - normal distribution

std - standard deviation

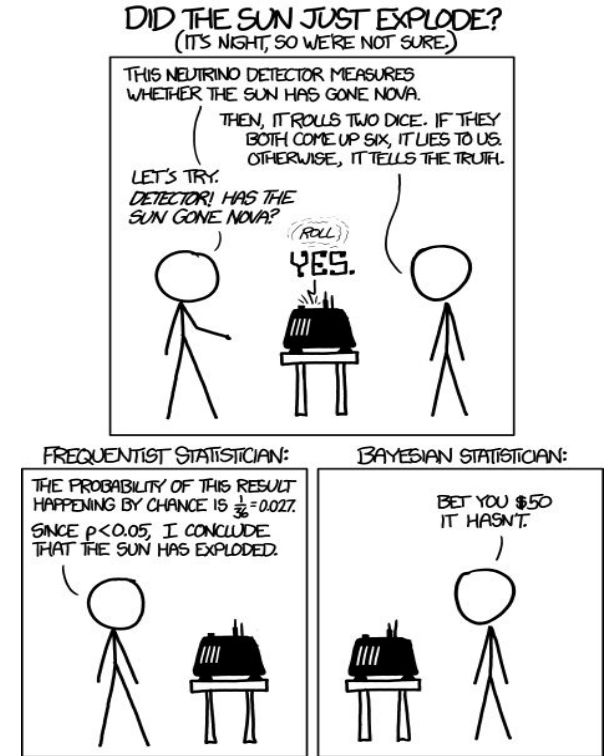
μ - mean

Frequentist

Repeatable random sample

Fixed data

$$y = mx + b$$



The Task

This is what you want your neural network to solve, the goal

Think carefully, you can get very different answers depending on how you frame your task.

-- First self driving cars were trained with human drivers and learned to drive on their own.

-- Test car is successfully driven down the road, then down a hill into a lake

? Car learned to follow the curb

Classification

Classification: image recognition, medical symptoms, marketing, voters

Features are turned into vectors and the distance between vectors is used to determine group membership. (city block, bird's eye)

Features representing values are scaled between 0,1 or -1,1

Features representing classes are converted to one hot vectors

Outputs a one hot vector selecting a class

Classification is unsupervised learning and might not always classify in a useful way

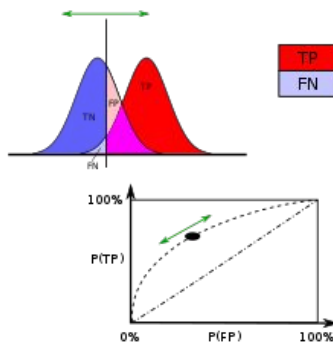
Different runs with data that's been shuffled between runs can classify the data differently

K-Means Clustering

pick n points at random (centroids), must give algorithm the number of clusters

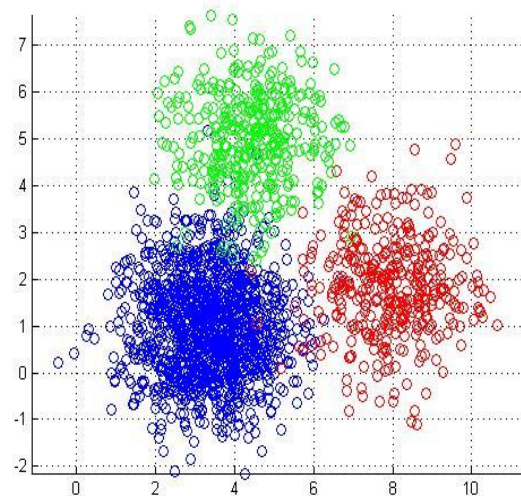
put each point in the group of the closest centroid using either Manhattan or Euclidean distance, repeat until stable

ROC Receiver operating characteristic ~ 1



<http://scikit-learn.org/stable/modules/generated/sklearn.cluster.KMeans.html>

<http://stanford.edu/~cpiech/cs221/handouts/kmeans.html>



Classification with missing inputs

Joint Training of Deep Boltzmann Machines for Classification, Goodfellow (Goodfellow 2013b paper)

<https://arxiv.org/pdf/1301.3568.pdf>

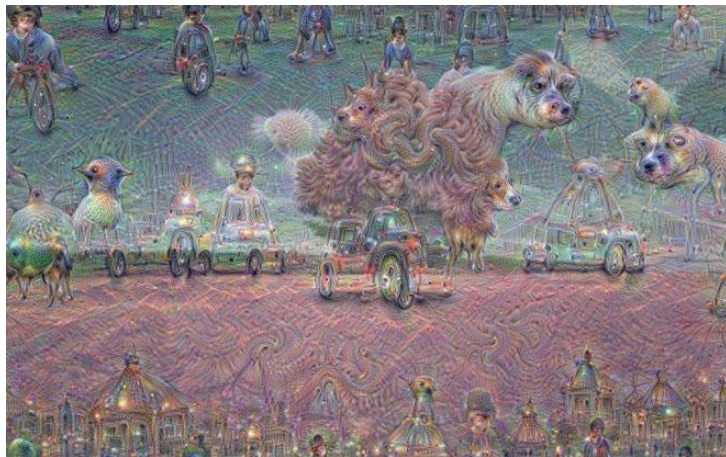
Restricted Deep Boltzmann Machines, Salakhutdinov and Hinton

<http://proceedings.mlr.press/v5/salakhutdinov09a/salakhutdinov09a.pdf>

Deep Learning 4j has a nice tutorial on Boltzmann machines

<https://deeplearning4j.org/restrictedboltzmannmachine>

(aka Deep Belief Networks)



Linear Least Squares Regression

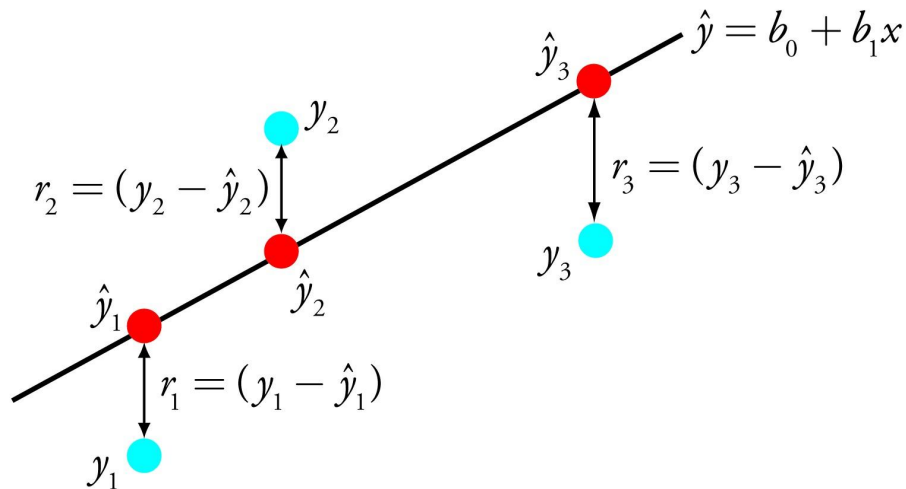
Solve for unknown $f(x)$

$$J(w) = \|y - Xw\|^2$$

Use derivative to find direction to change weights

$$dJ/dw = 2 * X^T X w - 2 * X^T y$$

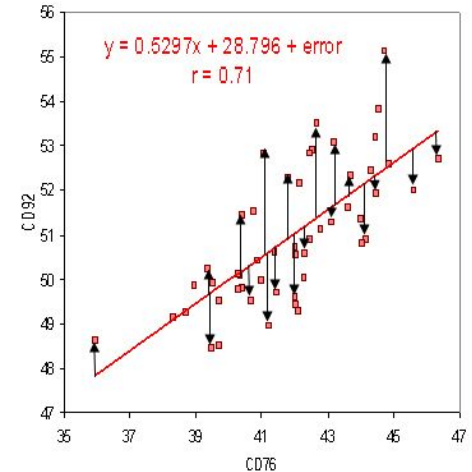
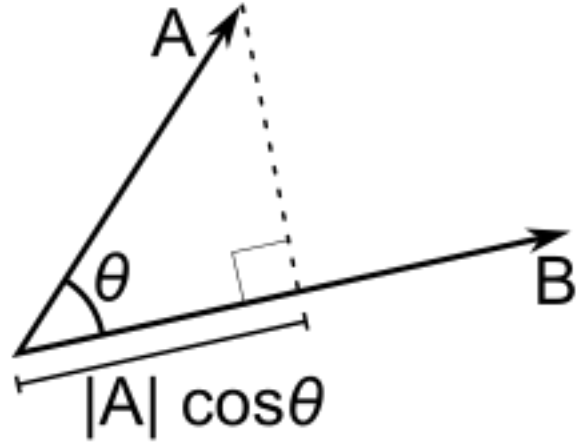
$$w = w - \text{learning_rate} * 2 * (X^T X w - X^T y)$$



Linear Regression

$$\hat{y} = \text{dot}(w.T, x)$$

$$\text{MSE} = 1/m * \text{Sum}(y_{\text{predicted}} - y_{\text{test}})^2$$



Transcription

OCR, Optical Character Recognition

Multi-digit Number Recognition from Street View

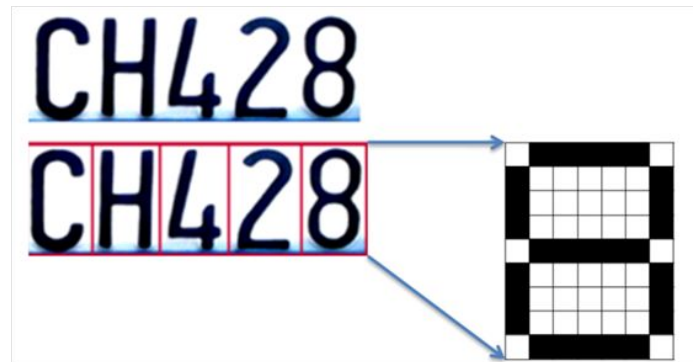
<https://static.googleusercontent.com/media/research.google.com/en//pubs/archive/42241.pdf>

Goodfellow talk on the paper

https://www.youtube.com/watch?v=vGPI_JvLoN0

Dataset

<http://ufldl.stanford.edu/housenumbers/>



Machine Translation

Used for natural language translation

Typically done with a recurrent network (RNN, LSTM, GRU, ...)

One language is the input, the second language is the output.

The network is trained back and forth rather than just forward feeding data and back feeding the errors

<https://nlp.stanford.edu/>

<https://github.com/MicrosoftTranslator/DocumentTranslator>

Structured Output

Collobert Home Page

<https://scholar.google.com/citations?user=32w7x1cAAAAJ&hl=en>

Natural Language Processing (Almost) from Scratch (Collobert 2012)

<http://www.jmlr.org/papers/volume12/collobert11a/collobert11a.pdf>

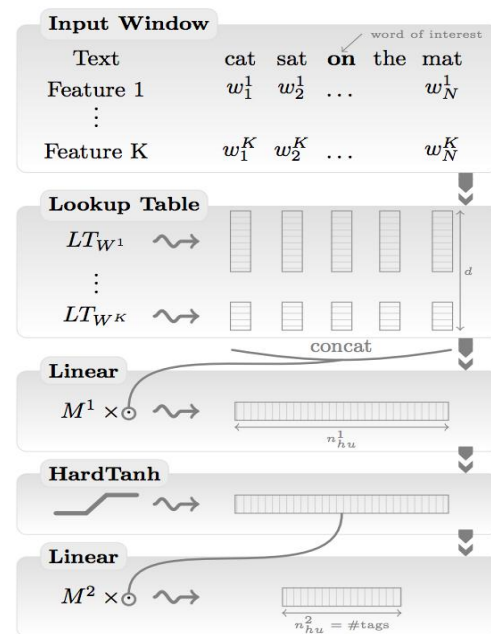


Figure 1: Window approach network.

Anomaly Detection

Make a prediction, points outside the prediction area are anomalies

RAD - Outlier Detection on Big Data

<https://medium.com/netflix-techblog/rad-outlier-detection-on-big-data-d6b0494371cc>

Bayesian Anomaly detection

<http://www.dtic.mil/cgi-bin/GetTRDoc?AD=ADA610860>

Synthesis and Sampling

Mostly used in the game industry to create landscapes, characters, etc Gamasutra

<https://www.gamasutra.com/topic/game-developer>

Microsoft Speech recognition and synthesis

<https://github.com/Microsoft/Windows-universal-samples/tree/master/Samples/SpeechRecognitionAndSynthesis>

<https://www.microsoft.com/en-us/research/blog/microsoft-researchers-achieve-new-conversational-speech-recognition-milestone/>

<https://arxiv.org/pdf/1708.06073.pdf>

Filling in missing values

- Use most common value, average or median
- Interpolation between values (time series)
- Drop that feature
- Infer from other features

The Performance Measure

~ 80% of data is used for training

~ 10% of data to validate

~ 10% of data is held out to be used as a test after training is done

While accuracy less than (pick a number)%:

- run training data in batches

- run validation every 100th epoch or so:

Then test on hold out data after training is complete

Training, testing, validation data must be statistically similar

The Experience

Unsupervised: data is not labeled

Supervised: labeled data aka target

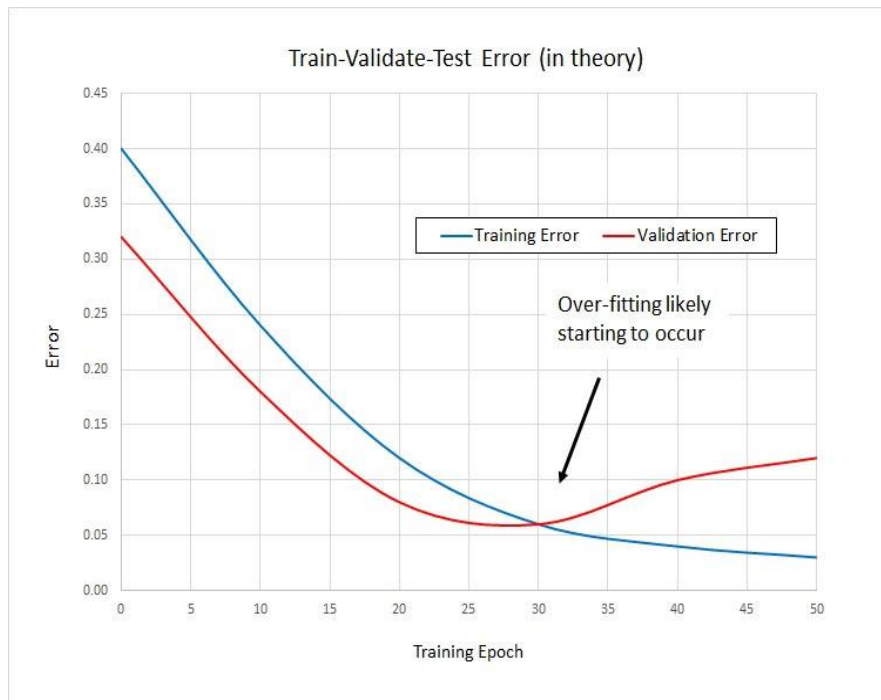
Chain rule of probability:

$$p(x) = \text{Product}(p(X_i | x_1, x_2 \dots))$$

Algorithm is learning the joint distribution of the data:

$$p(y|x) = p(x,y) / \text{Sum}(p(x,y))$$

Capacity, Overfitting, Underfitting



Confusion Matrix - very useful for debugging

| | | prediction outcome | | total |
|--------------|------|--------------------|----------------|-------|
| | | p | n | |
| actual value | p' | True Positive | False Negative | P' |
| | n' | False Positive | True Negative | N' |
| total | | P | N | |

No Free Lunch Theorem

Says that if you average over all the machine learning algorithms each will have the same performance on the same data set.

The only way one strategy can outperform another is by adjusting the network to the problem.

- convolutional networks for vision
- recursive networks for time series

<http://www.asimovinstitute.org/neural-network-zoo/>

Hyperparameters

Variables used to control the algorithm

- max, min, depth of branches on a decision tree
- number of layers, size of layers in network

Usually a grid search (nested loops over several options) is used with cross validation

http://scikit-learn.org/stable/modules/grid_search.html

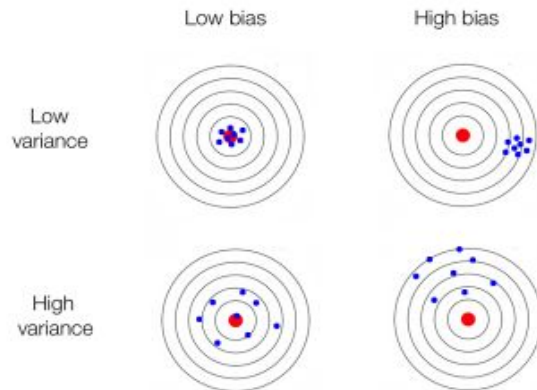
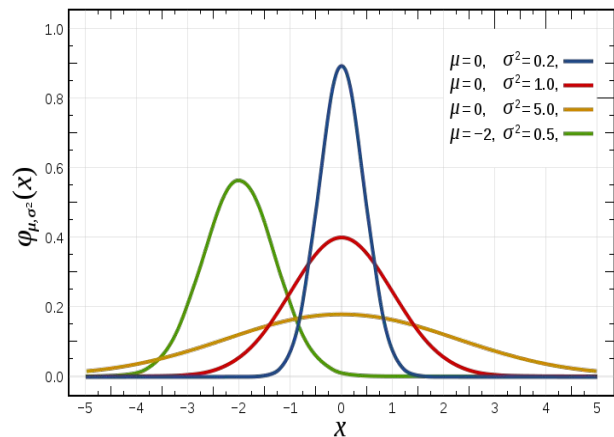
A recent development is that randomly picking values for a grid search is more effective than stepping through all values.

<http://www.jmlr.org/papers/volume13/bergstra12a/bergstra12a.pdf>

Bias vs Variance

High bias happens during under fitting, high variance during overfitting

Standard error is the square root of variance



Bias is hitting the 12 on the dart board every time, variance is hitting a wide area on the board

Entropy

Entropy = how many bits are needed to encode information

(* information theory uses Log2, not e, not 10)

$$\text{Entropy} = - \sum_i P_i \log_2 P_i$$

A Mathematical Theory of Communication, Shannon

<http://affect-reason-utility.com/1301/4/shannon1948.pdf>

KL Divergence

Kullback-Leibler Divergence tells how much information is lost between the neural network prediction and the actual values

$$D_{KL}(p||q) = \sum_{i=1}^N p(x_i) \cdot (\log p(x_i) - \log q(x_i))$$

The principle of maximum likelihood says that given the training data, we should use as our model the distribution $f(w)$ that gives the greatest possible probability to the training data

<https://www.countbayesie.com/blog/2017/5/9/kullback-leibler-divergence-explained>

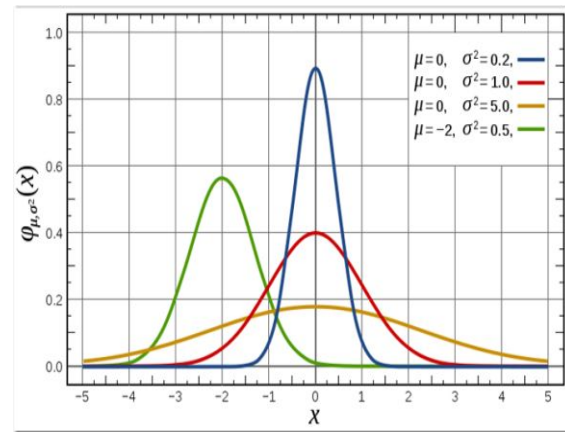
<http://cseweb.ucsd.edu/~elkan/250B/logreg.pdf>

Maximum Log Likelihood

$$\theta_{MLE} = \arg \max_{\theta} \log P(X|\theta)$$

$$= \arg \max_{\theta} \log \prod_i P(x_i|\theta)$$

$$= \arg \max_{\theta} \sum_i \log P(x_i|\theta)$$



Conditional Log Likelihood

Weights are trained to maximize the conditional data likelihood

$$W^* = \operatorname{argmax}_W \prod_{d \in D} P(Y^d \mid X_1^d \dots X_n^d)$$

which is the same as maximizing the conditional log likelihood

$$W^* = \operatorname{argmax}_W \sum_{d \in D} \ln P(Y^d \mid X_1^d \dots X_n^d)$$

<http://cseweb.ucsd.edu/~elkan/250B/logreg.pdf>

Mean Squared Error

If the data examples are i.i.d. ...

independent and identically distributed, then the Conditional Log Likelihood can be reduced to the Mean Squared Error

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \tilde{y}_i)^2$$

Independent Identically Distributed

Garbage in, garbage out

Independent

- features cannot be correlated

Identically Distributed

- training, validation, test data
- outputs (if 90% of the output is true, the network will always guess true and you'll have 90% accuracy - check the confusion matrix)

<https://pandas.pydata.org/pandas-docs/stable/generated/pandas.DataFrame.corr.html>

http://scikit-learn.org/stable/modules/generated/sklearn.metrics.confusion_matrix.html

Linear Regression as Maximum Likelihood

Dot product as a measure of similarity between vectors

MSE \sim Conditional log likelihood

As $n_{\text{samples}} \rightarrow \infty$ convergence increases

- True distribution of data must be in model family
- True distribution must correspond to one set of weights

Let's take the Con out of Econometrics

<http://www.econ.ucla.edu/workingpapers/wp239.pdf>

Bayesian Statistics

| The Posterior | The Evidence | The Prior |
|------------------------------------------------------------------------|-----------------------------------------------------------------------------|------------------------------------------------------------|
| | The probability of getting this evidence if this hypothesis were true | The probability of H being true, before gathering evidence |
| $P(H E)$ | $P(H E)$ | $P(H)$ |
| The probability that the hypothesis (H) is true given the evidence (E) | $P(E)$ | |
| | The marginal probability of the evidence (Prob of E over all possibilities) | |

Intuitive understanding of Bayes

https://arbital.com/p/bayes_frequency_diagram/?l=55z&pathId=28771

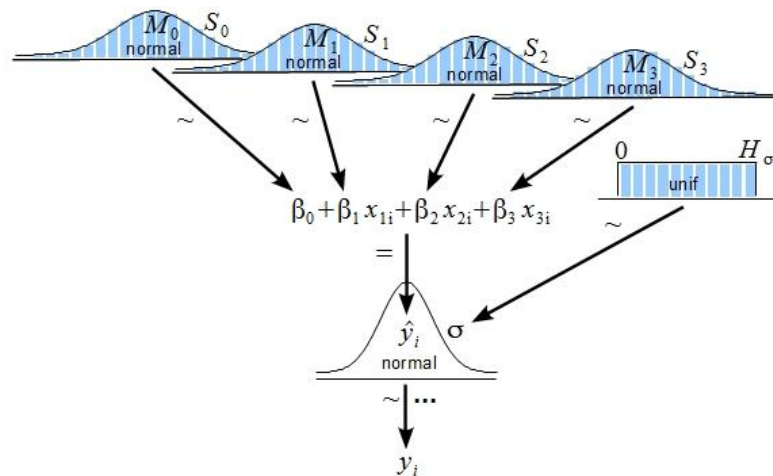
http://scikit-learn.org/stable/modules/naive_bayes.html

Bayesian Linear Regression

- Prior is uniform or Gaussian distribution with high entropy
- Observation reduces entropy and drives weights to single values

Using a known prior probability

Prior distribution shifts weights toward simpler smoother curves



<https://stats.stackexchange.com/questions/252577/bayes-regression-how-is-it-done-in-comparison-to-standard-regression>

Maximum a Posteriori (MAP) Estimation

$$\theta_{MAP} = \arg \max_{\theta} P(X|\theta)P(\theta)$$

$$= \arg \max_{\theta} \log P(X|\theta)P(\theta)$$

$$= \arg \max_{\theta} \log \prod_i P(x_i|\theta)P(\theta)$$

$$= \arg \max_{\theta} \sum_i \log P(x_i|\theta)P(\theta)$$

Maximum Likelihood and Maximum Posteriori Estimation

$$\begin{aligned}\theta_{MAP} &= \arg \max_{\theta} \sum_i \log P(x_i|\theta)P(\theta) \\ &= \arg \max_{\theta} \sum_i \log P(x_i|\theta) \text{ const} \\ &= \arg \max_{\theta} \sum_i \log P(x_i|\theta) \\ &= \theta_{MLE}\end{aligned}$$

<https://wiseodd.github.io/techblog/2017/01/01/mle-vs-map/>

Probabilistic Supervised Learning

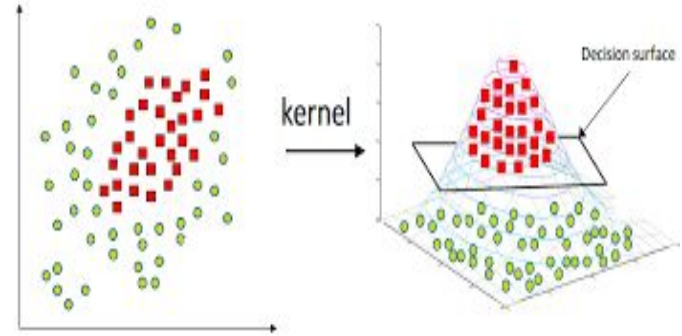
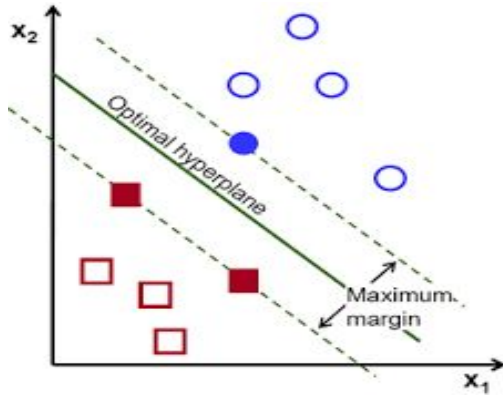
Normal (Gaussian Distribution)

Probability of y given x

The activation functions are used to create non-linearity

- sigmoid
- tanh
- relu
- softmax (used in last layer to scale values so they sum up to 1)

Support Vector Machines

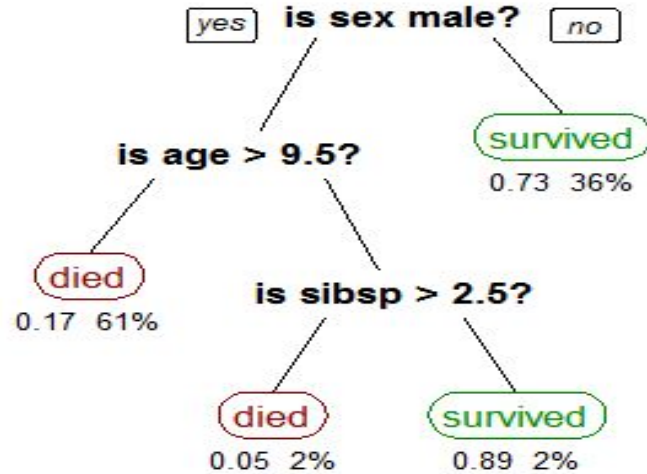


$O(n^2)$, can only separate 2 classes per SVM

<http://scikit-learn.org/stable/modules/svm.html>

https://docs.opencv.org/2.4.13.4/doc/tutorials/ml/introduction_to_svm/introduction_to_svm.html

Decision Tree



<http://scikit-learn.org/stable/modules/tree.html>

Using decision trees to see what a neural network is doing

<http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.12.6011&rep=rep1&type=pdf>

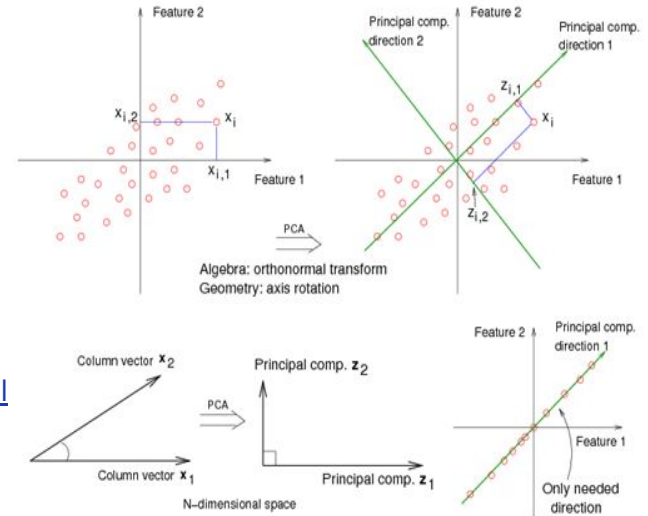
Principal Components Analysis

Used to compress data to reduce features by combining them in a way that maximizes variance

Use this when data is highly correlated to create independent feature representations

<https://onlinecourses.science.psu.edu/stat857/node/35>

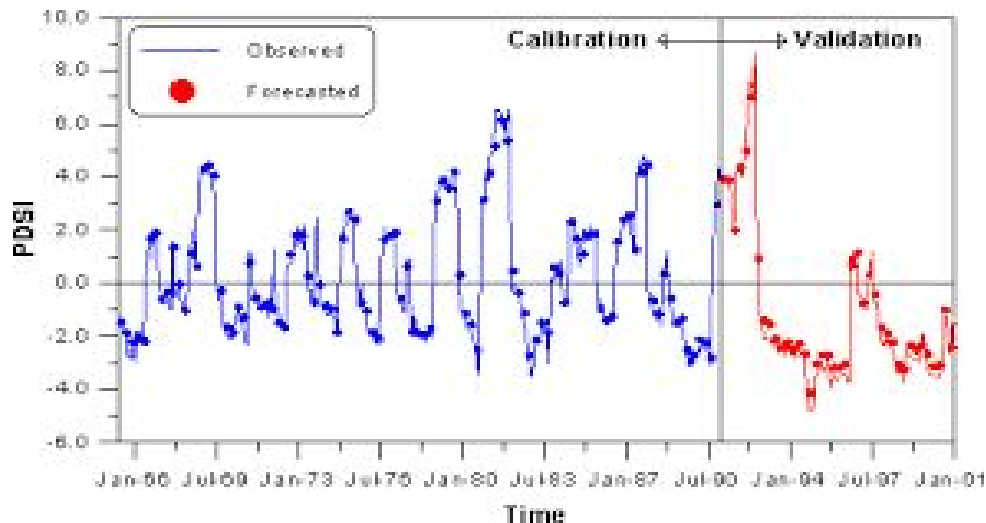
<http://scikit-learn.org/stable/modules/generated/sklearn.decomposition.PCA.html>



Time Series Predictions

Used to predict future values

- Must be stationary (rotate down to x axis)
- Generally use a log scale
- validation data must be future



Facebook's Prophet is a great open source tool for doing forecasts with large amounts of data

<https://facebook.github.io/prophet/>

Stochastic Gradient Descent

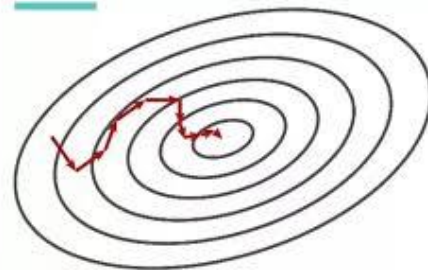
Repeat until convergence

{

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

}

Stochastic Gradient Descent



Building a Machine Learning Algorithm

Cleaned data

Multiply it by weights, iterative training for non-linear data

Minimize difference between cost function and actual data

(* no free lunch theorem)

New things being tried (or re-tried)

NEAT, NeuroEvolution of Augmenting Topologies <http://nn.cs.utexas.edu/downloads/papers/stanley.ec02.pdf>

Genetic Algorithms https://www.doc.ic.ac.uk/~nd/surprise_96/journal/vol1/hmw/article1.html

Reinforcement Learning <http://www.incompleteideas.net/book/bookdraft2018jan1.pdf>

Gradient Boosted Trees <http://xgboost.readthedocs.io/en/latest/model.html>

Curse of Dimensionality

This is the same issue you've seen in linear algebra.

If you have 2 equations and 3 unknowns there are an infinite number of solutions.

If you have a lot of features and few training samples there will be many paths through a network, they may validate and test okay, but will often randomly fail in real world use.

Regularization

$J(w)$ = Mean Standard Error + regularization

Training adjusts the weights to minimize the error, regularization forces the weights to stay small which forces the network to generalize

Used to prevent overfitting: To be useful a network must be able to generalize to new data.

- L1 regularization
- L2 regularization
- Early stopping
- Drop out
- Batch normalization

L1 Regularization

Minimize the sum of the absolute differences

Drives some weights to zero which has the effect of pruning features which acts as a feature selector. Almost never used in real problems

$$L1 = c1 * \text{sum}(\text{abs}(W))$$

L2 Regularization

Least squares: Minimize the square of the difference

Limits weight size without driving any to zero

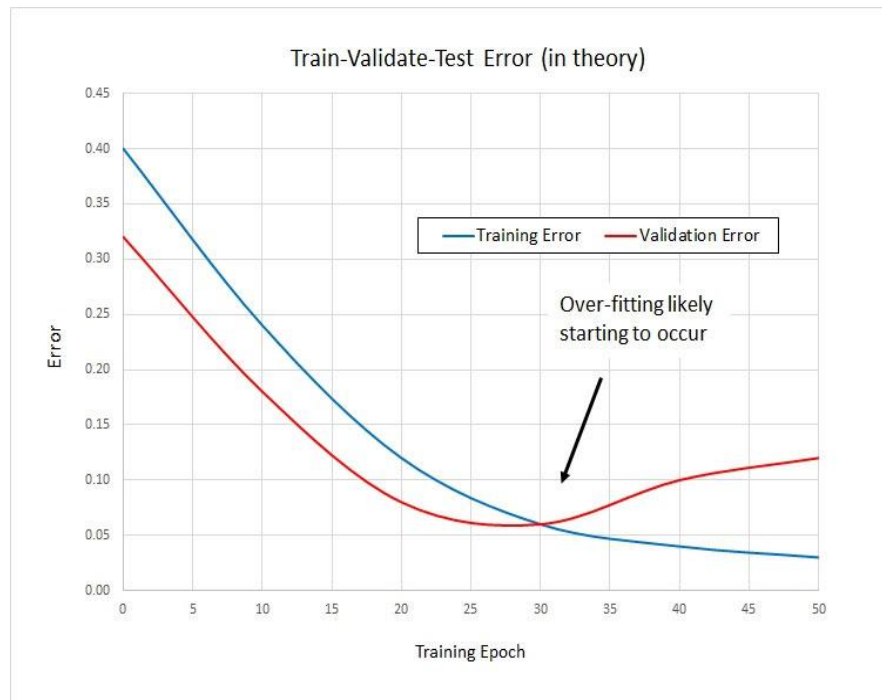
Rotational invariance is the main reason it is used instead of L1

$$L2 = c2 * \sum(W)^2$$

Early Stopping

Plot results during training

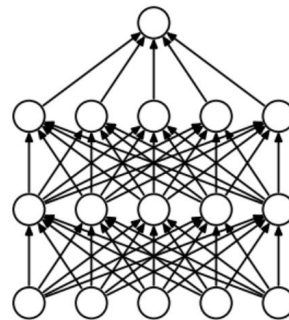
Stop training when training and validation errors converge



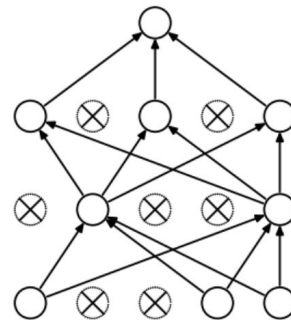
Dropout

Averages predictions over all weights

Randomly exclude a percent of hidden nodes from the neural network, change and randomly exclude a different group on next batch during training.



(a) Standard Neural Net



(b) After applying dropout.

Dropout: A Simple Way to Prevent NN from Overfitting

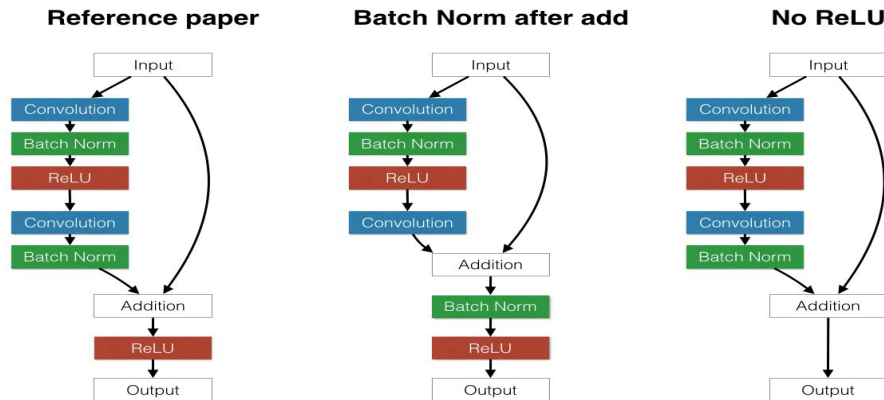
<https://www.cs.toronto.edu/~hinton/absps/JMLRdropout.pdf>

Batch Normalization

Is often used in place of DropOut, can use higher learning rates, networks train faster.

For each batch of training examples recenter the features.

In practice, the original input is re-fed to later layers



<https://arxiv.org/pdf/1502.03167.pdf>

Manifold Learning

Find a useful way to represent the data in less dimensions

Example: represent road as 2D in a 3D world

It's like PCA but attempts to retain more information

t-SNE Used often for visualization

Doesn't always work, information can be lost or distorted

<http://scikit-learn.org/stable/modules/manifold.html>