

AdaBoost

1. Assign every observation, x_i , an initial weight value, $w_i = \frac{1}{n}$, where n is the total number of observations.
2. Train a "weak" model. (most often a decision tree)
3. For each observation:
 - 3.1. If predicted incorrectly, w_i is increased
 - 3.2. If predicted correctly, w_i is decreased
4. Train a new weak model where observations with greater weights are given more priority.
5. Repeat steps 3 and 4 until observations are perfectly predicted or a preset number of trees are trained.

AIC

"Ey-Kye-Ih-Key"

Information Criteria

$$AIC = \frac{1}{n\hat{\sigma}^2} (RSS + 2d\hat{\sigma}^2)$$

Residual Sum of Squares
number of features
sample variance
number of observations

Used to compare which model is better. For example during feature selection.

ARCHITECTURE OF A NEURAL NETWORK

- The architecture of a neural network refers to the units, their activation functions, how many layers etc.
- Most neural networks' architecture can be understood of as stacked layers of units.
- The best architecture for a problem should be found through experimentation using validation sets.

Chris Albon

BAGGING vs. DROPOUT

MODELS: In bagging, all models are independent.

In dropout, subnetworks share parameters.

TRAINING: In bagging, all models are trained.

In dropout, only a fraction of possible subnetworks are trained.

BAG OF WORDS

Converts text to a matrix where every row is an observation and every feature is a unique word. The value of each element in the matrix is either a binary indicator marking the presence of that word or an integer of the number of times that word appears.

Chris Albon

BIAS - Variance Tradeoff

$$\text{Error}(x) = \left(E[\hat{f}(x)] - f(x) \right)^2 + E\left[\hat{f}(x) - E[\hat{f}(x)] \right]^2 + \sigma_e^2$$

↓ ↓ ↓ ↓
predicted true predicted average predicted value irreducible error

Bias²

How much predicted values differ from true values.

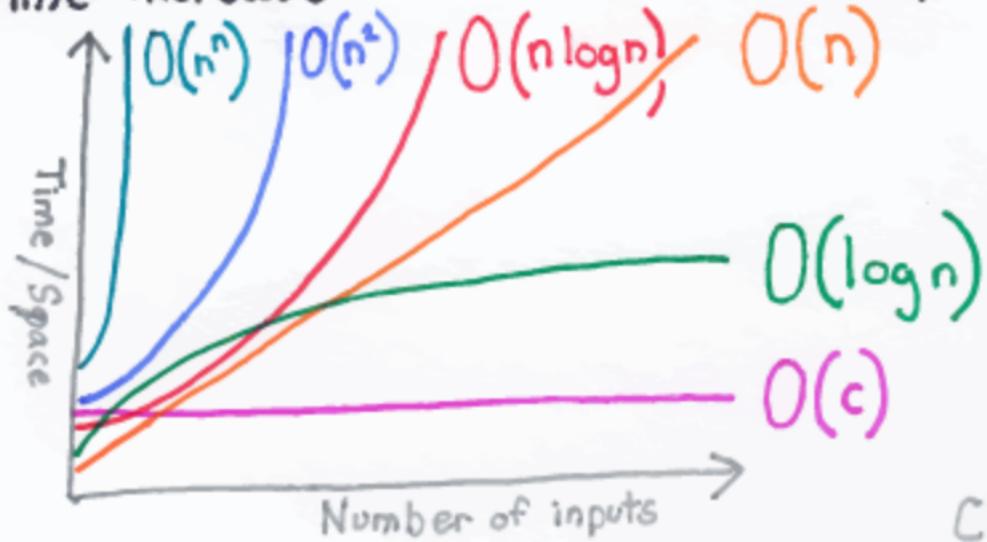
Variance

How predictions made on the same value vary on different realizations of the model

BIG O



Representation of an algorithm's space requirement or run time increase as the number of inputs increase



ChrisAlbon

BRIER SCORE

$$BS = \frac{1}{n} \sum_{t=1}^n \left(P_t - O_t \right)^2$$

Predicted Probability
Actual outcome

↑
number of observations

Brier score shows the squared mean difference between the predicted probability of all observations with their actual outcome. The lower score the better. Ranges between 0 and 1.

Chris Albon

COMBINATION

Number of subsets of k items from a total of n items.

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

"From n , choose k ."

CONDITIONING

Conditioning is a measure of how much a function's outputs change when its inputs change. Poorly conditioned functions are highly sensitive to rounding errors that can happen in numerical computing.

ChrisAlbon

COST AND LOSS FUNCTIONS

Cost functions and loss functions mean the same thing. They are the objective function we are trying to train an model that minimizes.

EXAMPLE: Cross-entropy

Chris Albon

Ep

C_p is used in model selection to compare the performance of different models.

$$C_p = \frac{1}{n} \left(RSS + 2d\hat{\sigma}^2 \right)$$

The penalty to adjust for the fact that training data underestimates the test error.

Number of observations

Residual Sum of squares

Number of features

Estimated error variance

CDF

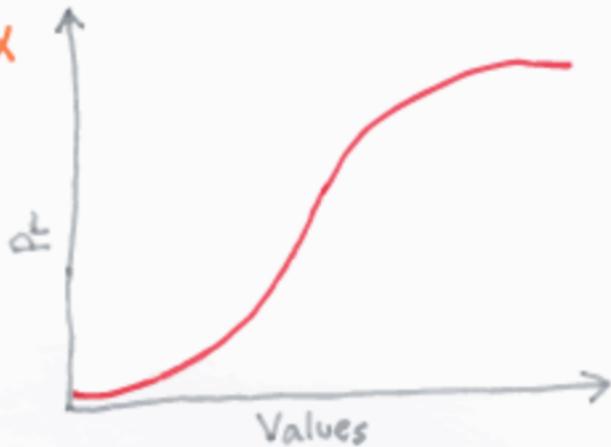
Cumulative Distribution Function

CDF tells us the probability a random variable returns a value less than some specified value. It is the accumulation of the probability of values up to some value.

$$f(x) = \Pr[X \leq y]$$

↑
Random variable

Some value $\in X$



DBSCAN

DBSCAN looks for densely packed observations and makes no assumptions about the number or shape of clusters.

1. A random observation, x_i , is selected
2. If x_i has a minimum of close neighbors, we consider it part of a cluster.
3. Step 2 is repeated recursively for all of x_i 's neighbors, then neighbors' neighbors etc... These are the cluster's core members.
4. Once Step 3 runs out of observations, a new random point is chosen

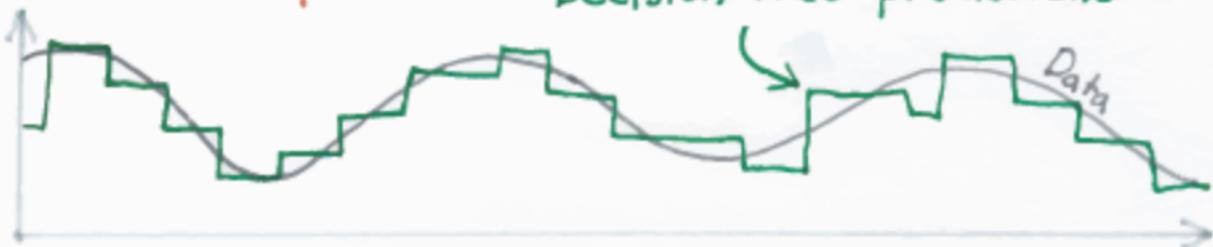
Afterwards, observations not part of a core are assigned to a nearby cluster or marked as outliers.

ChrisAlbon

DECISION TREE REGRESSION

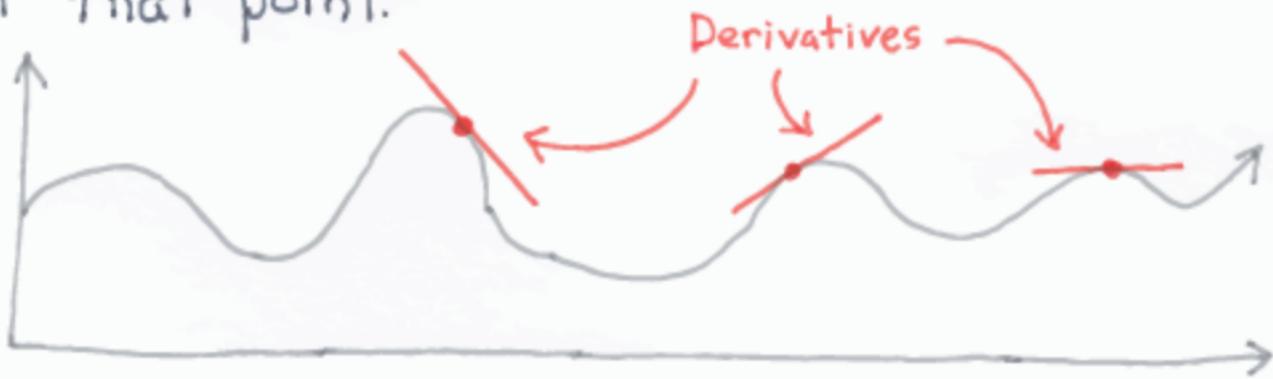
Similar to decision tree classification, however uses Mean Squared Error or similar metrics instead of cross-entropy or Gini impurity to determine splits.

Decision tree predictions



DERIVATIVE

The derivative of a function is its rate of change. Intuitively, the first derivative at some point is the slope of the function at that point.



DETERMINANTE

Determinants can be thought of as a scalar value summarizing a matrix. Formally, the determinant of a matrix is the product of all the matrix's eigen-values.

The absolute value of the determinant can be imagined as how much the matrix expands or contracts space.

DEEES K-NN LEARN

k-nearest neighbor does not "learn" per-se. It is lazy and just memorizes the data.

Chris Albon

EARLY STOPPING

ADVANTAGES

1. Does not require altering the network architecture or training method.
2. Stops automatically at the best point rather than requiring hyperparameter tuning like weight decay.

ChrisAlbon

THE EFFECT OF ONE-HOT ON

FEATURE IMPORTANCE

In random forests we can measure the importance of each feature. However, if our feature is a nominal categorical feature that we one-hotted, the importance of that feature will be distributed across all of these one-hotted features.

ChrisAlbon

ELASTICNET

A linear regression model that combines the L1 and L2 regularizers.

$$\text{RSS} + \underset{\substack{\downarrow \\ \text{Alpha}}}{{\alpha}} \underset{\substack{\downarrow \\ \text{Rho}}}{{P}} \|\underset{\substack{\uparrow \\ \text{Weights}}}{{w}}\|_1 + \frac{\alpha(1-P)}{2} \|w\|_2^2$$

↑
Residual sum of squares

Alpha determines the regularization strength.

Rho determines the balance between L1 and L2.

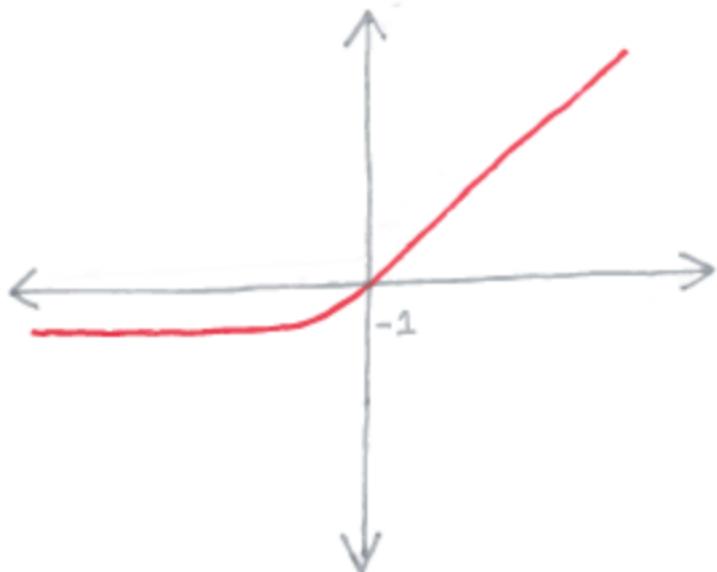
ChrisAlbon

ELU's

Exponential Linear Units

$$\phi(z) = \begin{cases} z & \text{if } z \geq 0 \\ \alpha [\exp(z) - 1] & \text{otherwise} \end{cases}$$

positive hyperparameter



EPOCH

In artificial neural networks each time all observations have been sent through the network is called an epoch. Training neural networks typically involves multiple epochs.

ChrisAlbon

EXPLAINED SUM OF SQUARES

ESS measures the amount of variance (information) in the model

$$ESS = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

↑
Predicted value ↑
 mean value

EXPLAINED SUM OF SQUARES

ESS measures the amount of variance (information) in the model

$$ESS = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

↑
Predicted value ↑
 mean value

FPR

False Positive Rate

$$FPR = \frac{\text{False Positives}}{\text{False Positives} + \text{True Negatives}}$$

FEATURE SELECTION STRATEGIES

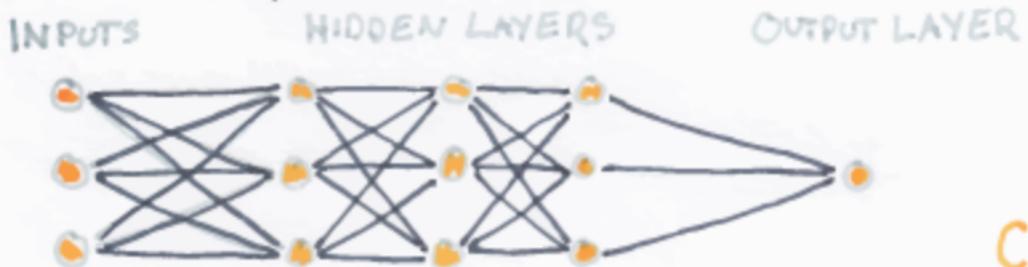
1. Remove highly correlated variables.
2. Run OLS and select significant features.
3. Forward selection and backwards selection.

↓
or recursive
4. Random Forest feature importance.
5. Lasso.

BY CHRIS ALBON

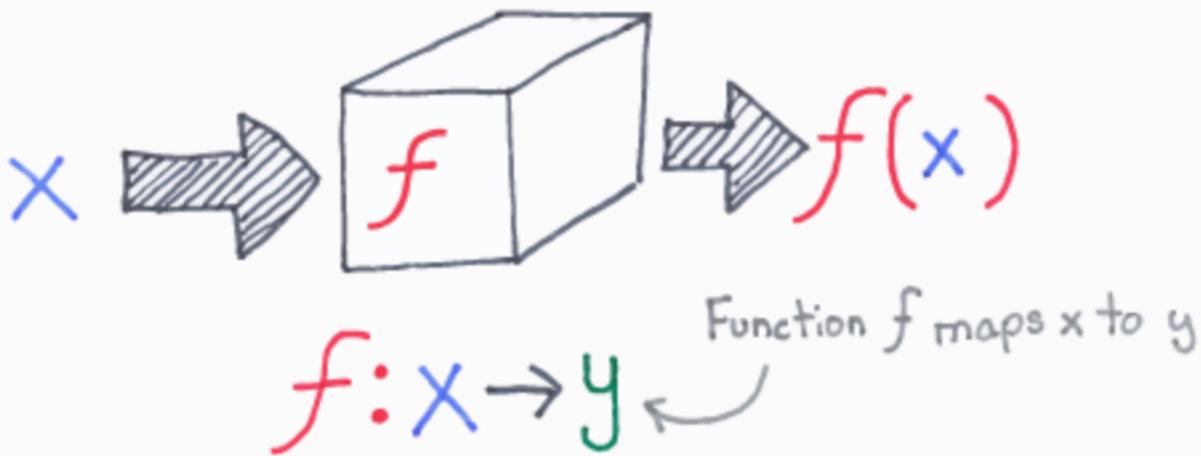
FEEDFORWARD NEURAL NETWORKS

- Also called multilayer perceptron.
- Called feedforward because information moves forward from inputs to output layer.
- Structure example:



Chris Albon

FUNCTION



- The domain of the function is the possible set of values of x .
- The range of the function is the possible set of values of y .

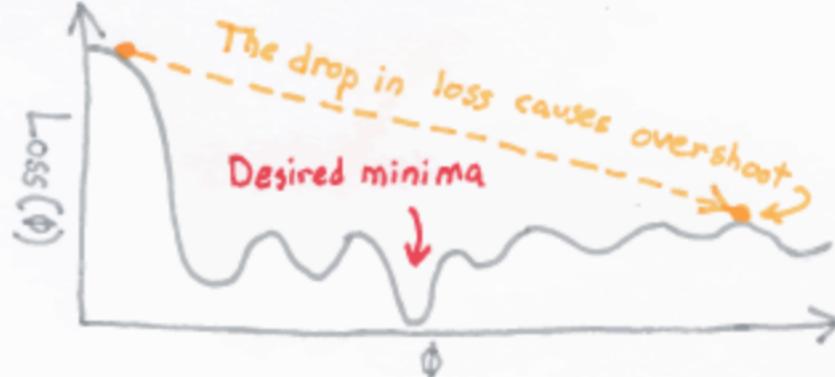
GRADIENT

The gradient of a function, f ,
is a vector containing all its
partial derivatives at some point, x .

$$\nabla_x f(x)$$

GRADIENT

When there is a steep drop in the loss function. Can cause a problem when the steep gradient causes the optimization algorithm to overshoot the minima.



Chris Albon

GRADIENT CLIPPING

Cliffs in the loss function can cause the learning process to overshoot a desirable minimum. This happens because the gradient at the cliff is so high. One solution is to prevent the gradient from taking extreme values:

$$\text{if } \|g\| > v: \quad g \leftarrow \frac{g}{\|g\|} \cdot v$$

gradient
Threshold
norm of gradient

ChrisAlbon

GREEK 3

LETTERS

N ν nu

Π π pi

Ξ ξ xi

ρ ρ rho

ο ο omicron

Σ σ sigma

HAMMING LOSS

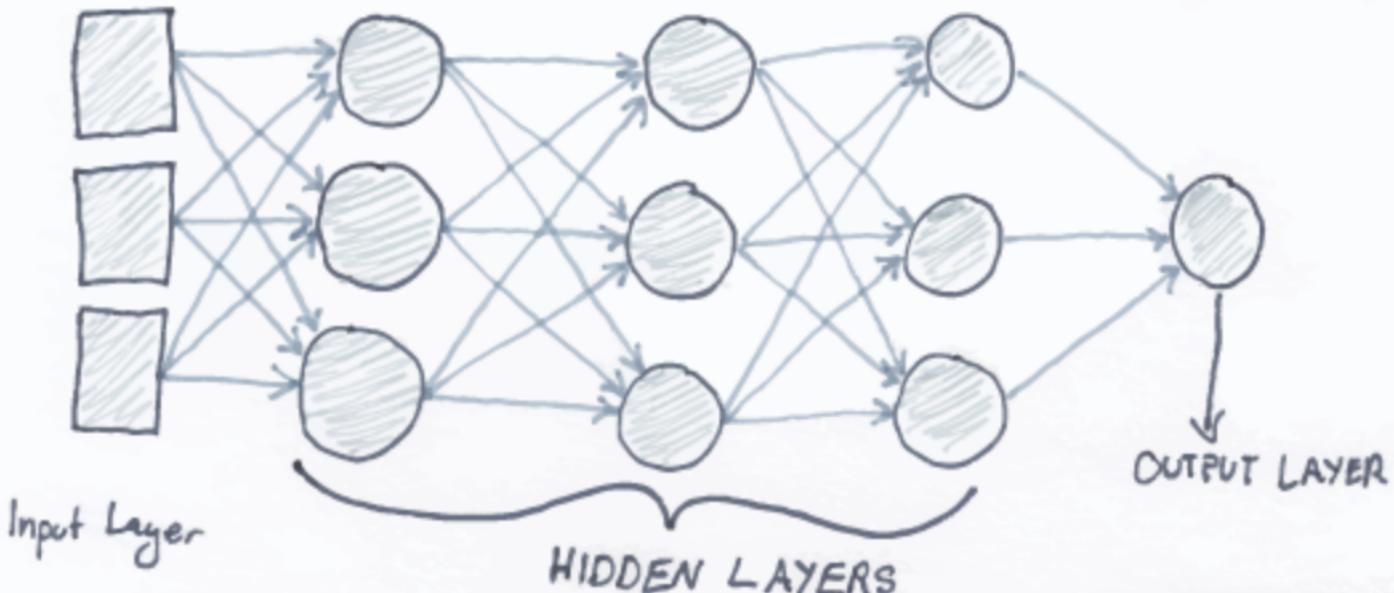
$$L_{\text{Hamming}} = \sum_{i=1}^n I(\hat{y}_i \neq y_i)$$

Predicted class true class

n ←
number of observations

HIDDEN LAYER

VALUES
CALLED HIDDEN LAYERS BECAUSE^A NOT IN THE DATA. RATHER,
VALUES ARE DETERMINED BY WHAT IS USEFUL FOR MODELING RELATIONSHIP.



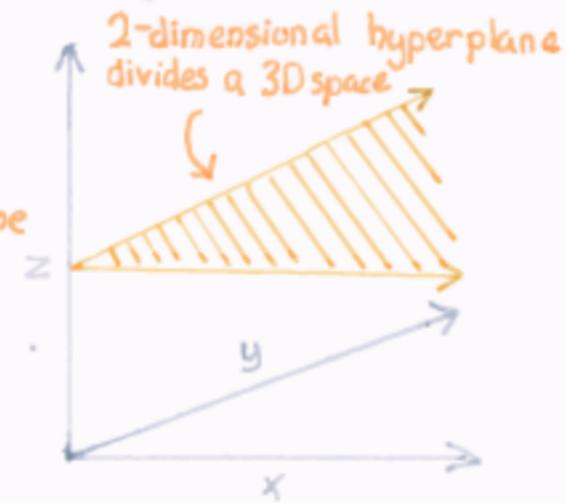
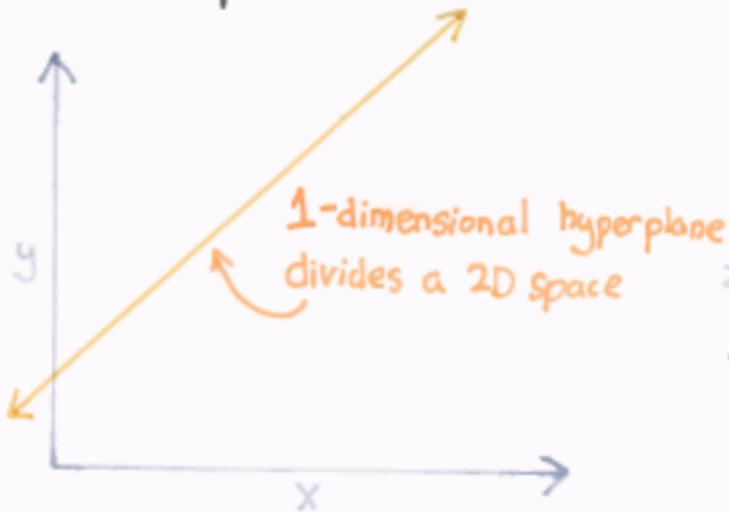
HOW NORM PENALTIES WORK

- L1 and L2 norm penalties shrink parameters toward zero.
- The benefit comes from less variance in parameter values, not necessarily small values.
- Zero is typically used because it is neither positive or negative.

Chris Albon

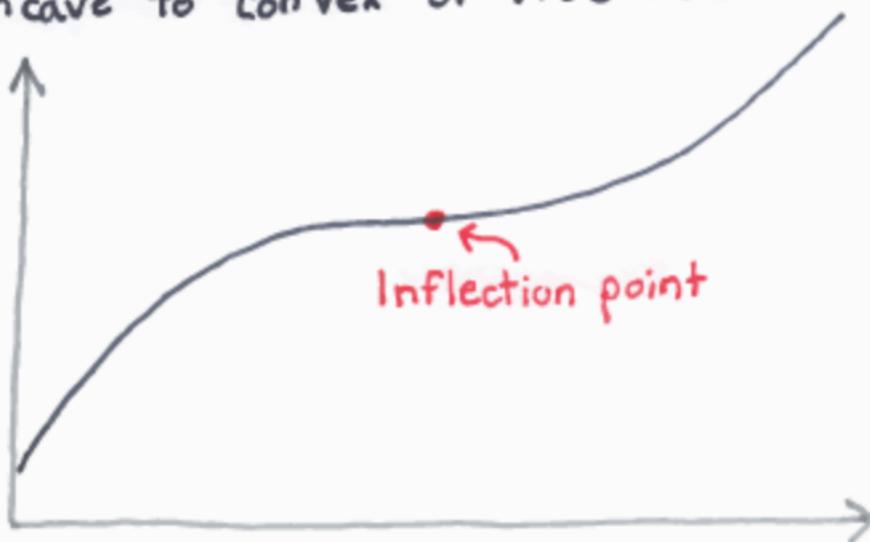
HYPERPLANE

In an n -dimensional space, a hyperplane is an $n-1$ plane that divides that space.



INFLECTION POINT

The point on a function where the surface changes from concave to convex or vice versa.



Chris Albon

INITIALIZING WEIGHTS IN FEEDFORWARD NEURAL NETWORKS

- Initialize with small random numbers.
- Common to draw initial weights from normal distribution.
- Biases initialized as zero or small positive numbers.

Chris Albon

INTERACTION TERM

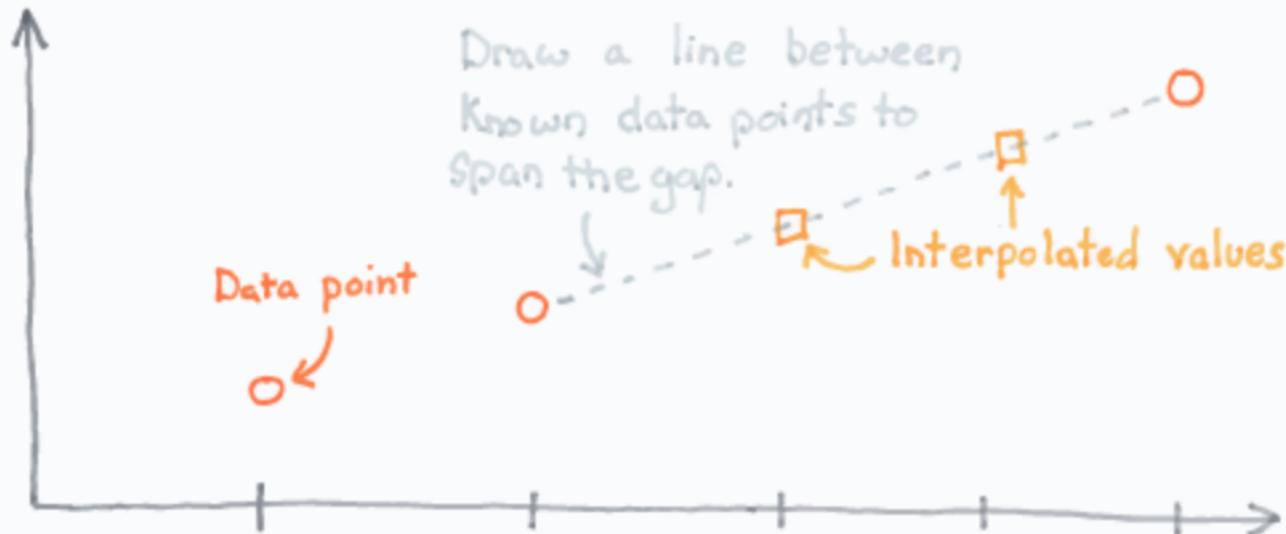
Interaction terms allow us model relationships when the effects of a feature on the target is influenced by another feature.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + e$$

The interaction of features
 x_1 and x_2 .

INTERPOLATION

A strategy to fill gaps of missing values by drawing a line between known values.



K-FOLD

CROSS-VALIDATION

$$CV_{(K)} = \frac{1}{K} \sum_{i=1}^K \text{Loss}_i$$

of folds
10 folds is common

For example:

- Mean squared error
- Log-loss
- Accuracy

K-NEAREST NEIGHBORS

TIPS AND TRICKS

1. All features should use the same scale.
2. K should be odd to avoid ties.
3. Votes can be weighted by the distance to the neighbor so closer observations' votes are worth more.
4. Try a variety of distance measurements.

Chris Albon

KNN

NEIGHBORHOOD SIZE

Small



$K =$ Low Bias, High Variance



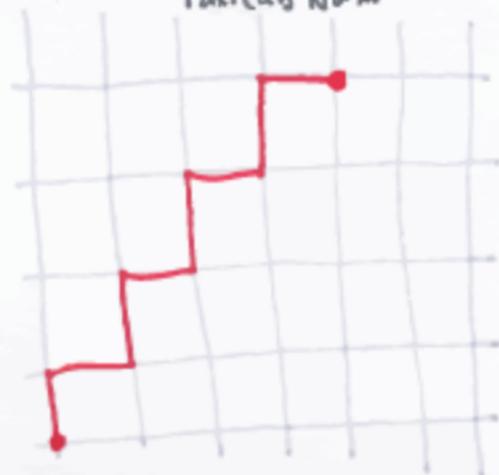
$K =$ High Bias, Low Variance

L₁Norm

(Manhattan Norm)

Also Called
"Taxicab Norm"

$$\|x\|_1 = \sum_{i=1}^n |x_i|$$



ChrisAlbon

L2 NORM

(Euclidean Norm)

$$\|x\|_2 = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2}$$

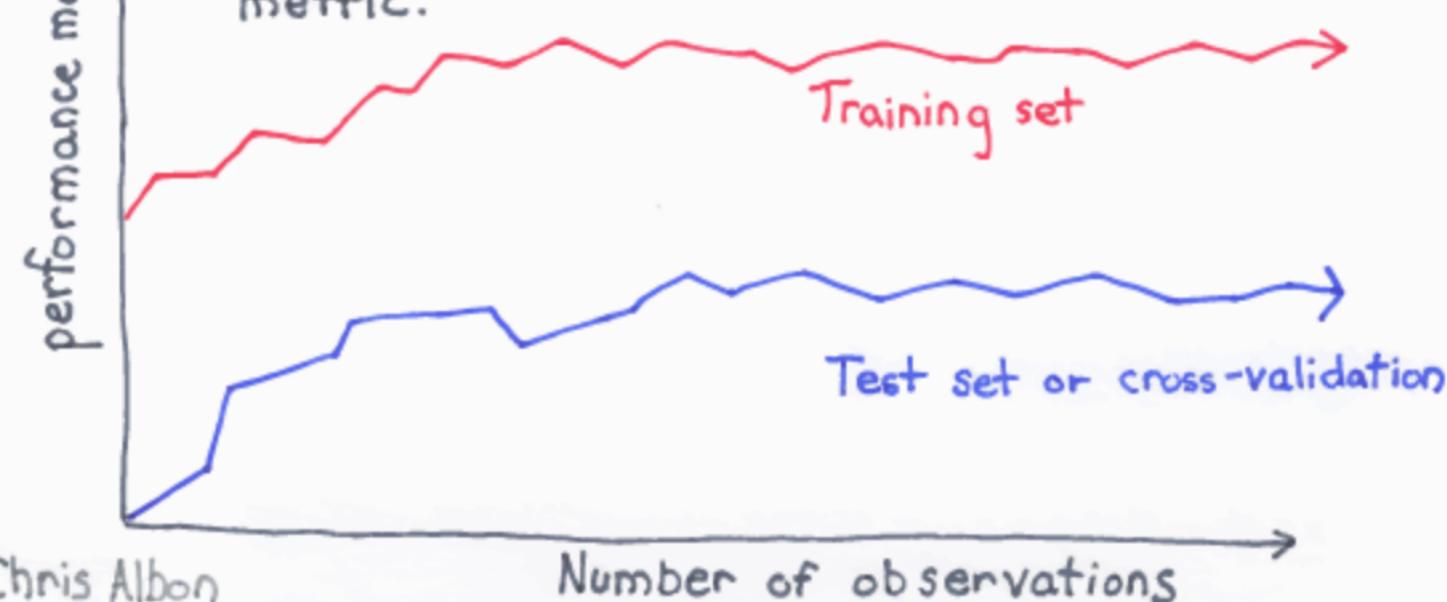
L2 norm is used in many places in machine learning such as normalizing observations and regularization.

Like in NLP.

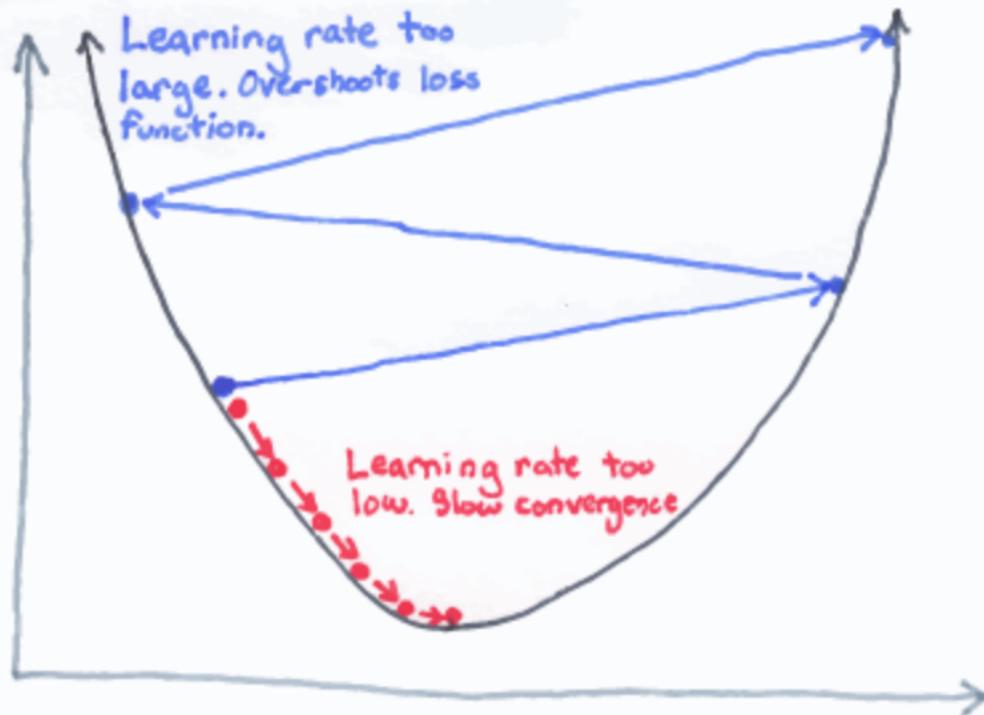
Example: Ridge Regression.

LEARNING CURVE

Learning curves visualize the effect of the number of observations on the performance metric.

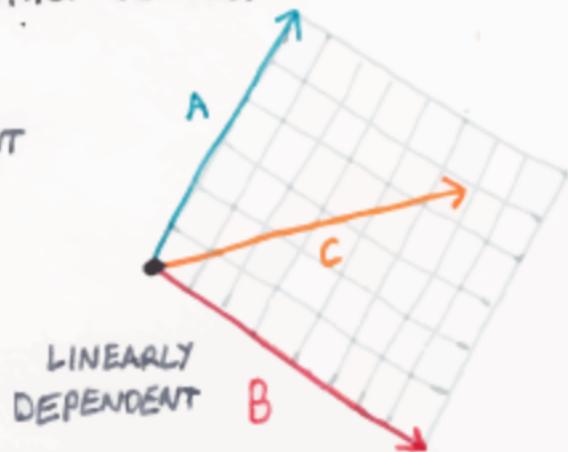
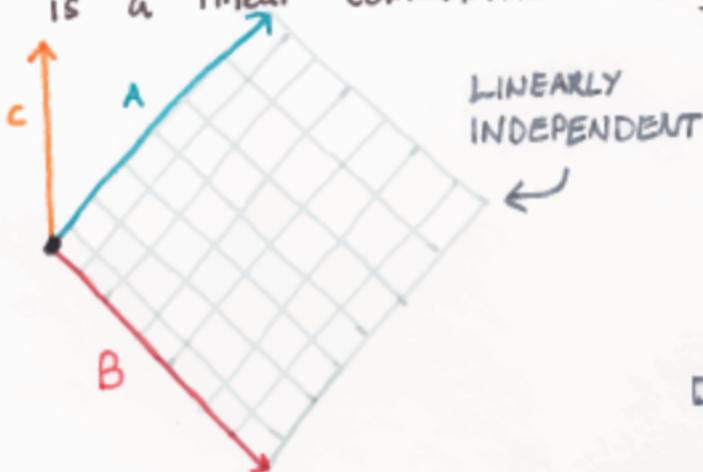


LEARNING RATE



LINEARLY INDEPENDENT

A matrix is linearly independent if no vector in the matrix is a linear combination of other vectors.



LINEAR COMBINATION OF A SET OF VECTORS

$$\sum_i c_i v^{(i)}$$

Scalar

Set of vectors

ChrisAlbon

MATRICES

Two dimensional array of scalars.

$$A = \begin{bmatrix} a_{1,1} & \dots & a_{1,n} \\ \vdots & \ddots & \vdots \\ a_{m,1} & \dots & a_{m,n} \end{bmatrix}$$

rows columns

Scalar

MATRIX INVERSE

$$A^{-1} A = I_n$$

Inverse of matrix A Matrix n-dimensional identity matrix

ChrisAlbon

MATRIX

MULTIPLICATION

$$C_{i,j} = \sum_k A_{i,k} B_{k,j}$$

row index column index

Matrix Matrix

Matrix multiplication is associative but not commutative.

$$A(BC) = (AB)C \quad AB \neq BA$$

MATTHEWS CORRELATION COEFFICIENT

$$M = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

↑ ↑ ↑ ↑
True Positive True Negative False Positive False Negative

Ranges from -1 (perfectly wrong classifier) to 1 (perfectly right classifier).

MEANSHIFT CLUSTERING

BY ANALOGY

I imagine a foggy football field with 100 people standing on it. Because of the fog, people can only see a short distance. Every minute each person looks around and takes a step in the direction of the most people they can see. As time goes on, people start to group up as they repeatedly take steps towards larger and large crowds. The end result is clusters of people around the field.

ChrisAlbon

MEAN ABSOLUTE

ERROR

$$\text{MAE}(y, \hat{y}) = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

↑ ↑ ↑
Number of Observations True Target Value Predicted Target Value

MSE

Mean

Squared
true y

Error

Predicted y
2

Squared

$$MSE = \frac{1}{n} \sum (y_i - \hat{f}(x_i))^2$$

number of observations

Also called \hat{y}_i
"y hat"

Common evaluation metric in regression

MIN MAX SCALING

Rescales feature values to between 0 and 1

$$x'_i = \frac{x_i - \min(x)}{\max(x) - \min(x)}$$

Original value
Minimum value in feature
Rescaled value
Maximum value in feature

MNAR

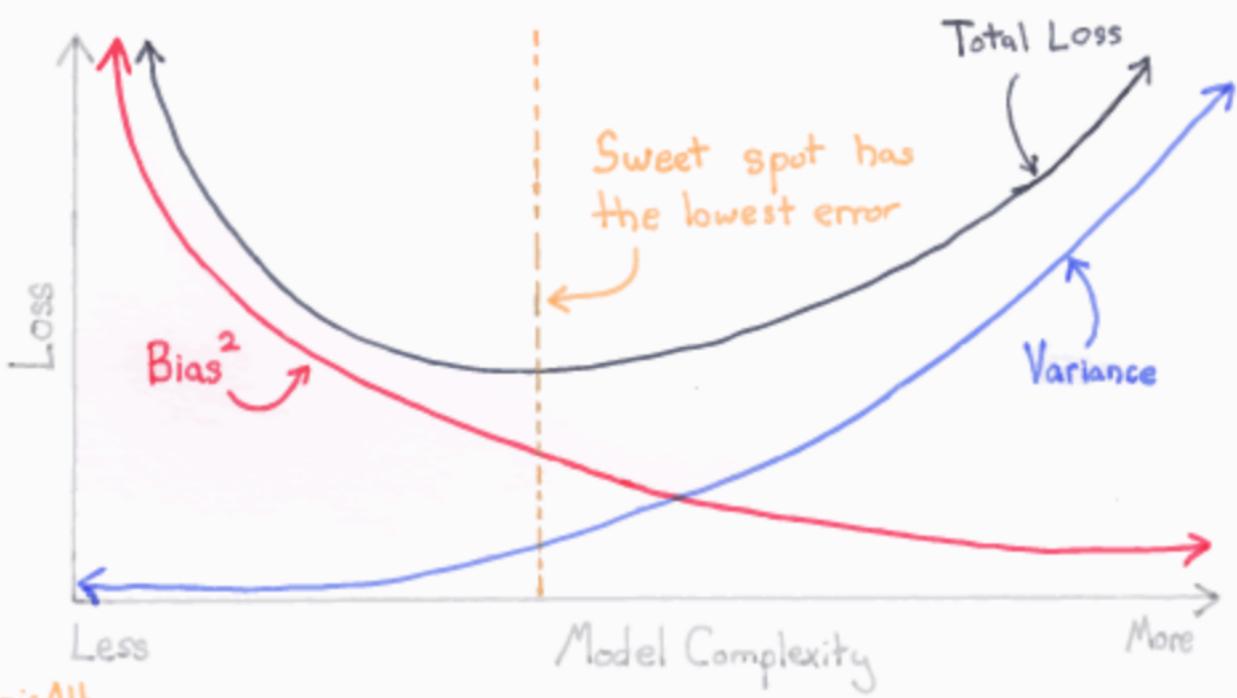
MISSING NOT AT RANDOM

A type of missing data where the probability a value is missing is not random and depends on information not captured in the other features.

For example, men are less likely to answer a salary question in a survey but we do not capture gender identity in another feature.

Chris Albon

MODEL COMPLEXITY



MODEL SELECTION

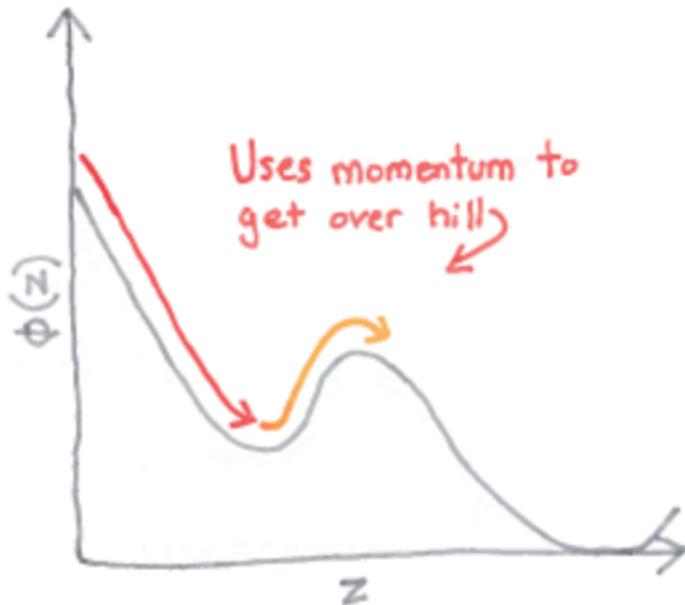
Finding the machine learning algorithm and its hyperparameter values that produce the best model.

ChrisAlbon

MOMENTUM

Momentum uses previous gradients to influence the movement of the optimizer.

Often used with stochastic gradient descent.



MOTIVATION FOR DEEP LEARNING

“Shallow learning” algorithms like support vector machines work well when we have well structured feature data. However, they often perform poorly in high dimensional spaces such as those found in computer vision and machine translation.

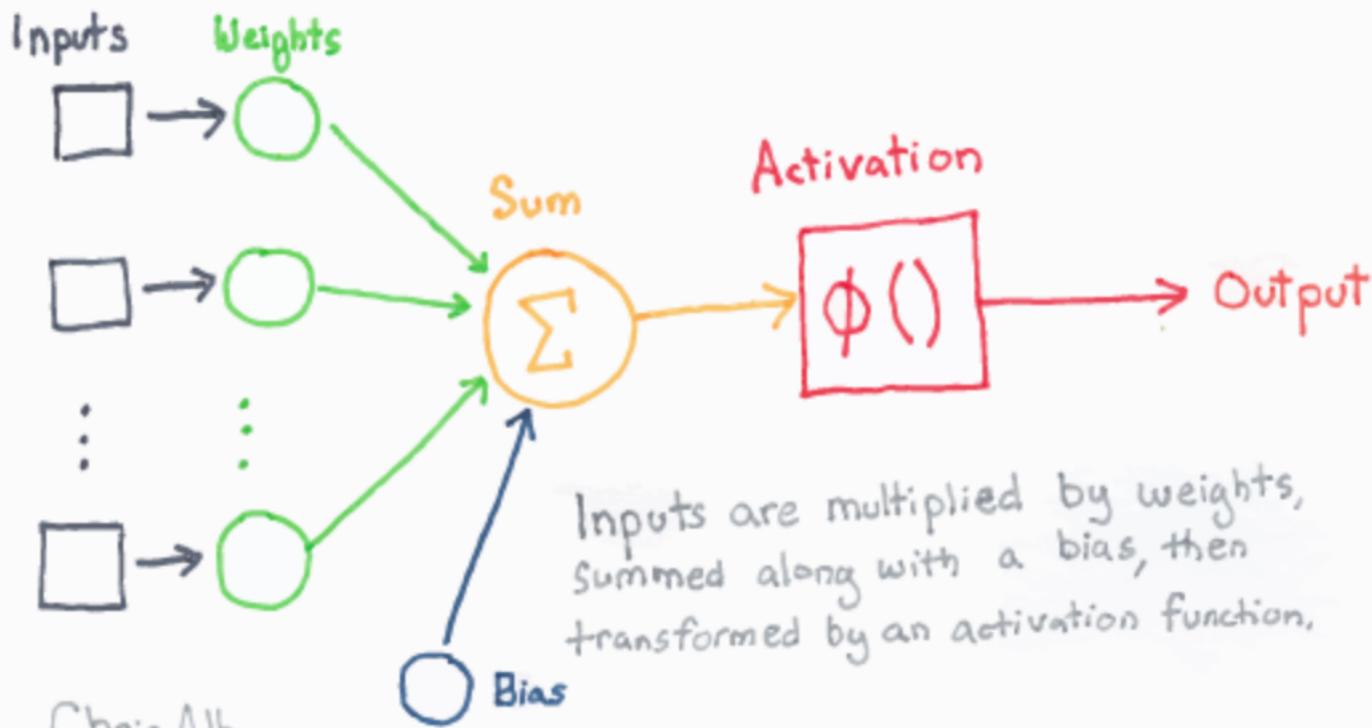
Chris Albon

NATURAL LOG

$$\ln(e^x) = x$$

Euler's number: 2.71828...

NEURON



NOISY RELU

Similar to ReLU but adds Gaussian noise:

$$\phi(z) = \max(0, z + N)$$

Value drawn from
a normal distribution:

$$N(0, \sigma(z))$$

NORMALIZED INITIALIZATION OF NEURAL NET PARAMETERS

Uniform distribution

$$W_{i,j} \sim U\left(-\sqrt{\frac{6}{m+n}}, \sqrt{\frac{6}{m+n}}\right)$$

Parameter in
a fully-connected
layer.

Number of inputs

Number of outputs

Chris Albon

NORMALIZING OBSERVATIONS

Rescaling the feature values of each observation so that they have a unit norm.

Two common norm values are L1 and L2.

$$\|x\|_2 = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2}$$

↑
L2 norm

ODDS RATIO

$$\frac{\Pr(x_1)/\Pr(\sim x_1)}{\Pr(x_2)/\Pr(\sim x_2)}$$

Odds of
event x_1

Odds of
event x_2

ChrisAlbon

OUT-OF-CORE

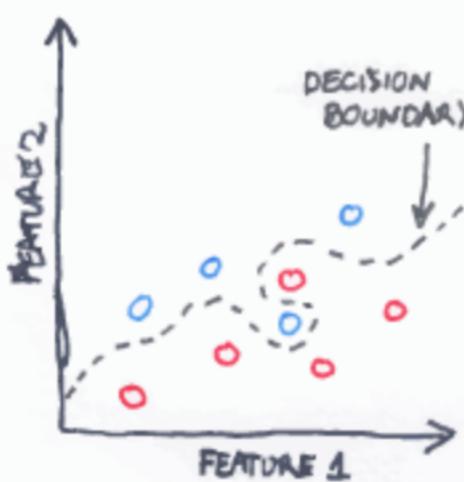
Learning methods when the data is too big to fit in a computer's RAM.

Examples of out-of-core options:

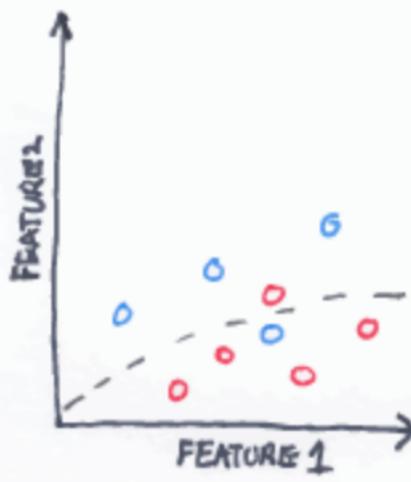
- Preprocess data in chunks
- Read and preprocess data line by line
- Incremental learning
- Stochastic techniques
- Partial fit learning methods

Chris Albon

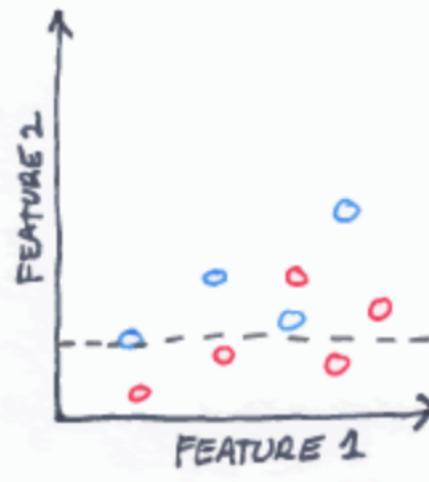
OVERFIT vs UNDERFIT



OVERFIT
"HIGH VARIANCE"



IDEAL



UNDERFIT
"HIGH BIAS"

CORRELATION

(Pearson's R)

$$\text{Cor}(X, Y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}}$$

value of x_i mean of x value of y_i mean of y

Ranges between -1.0 and +1.0. The closer to 0.0 the less linear dependence between variables.

PRECISION

Precision is the ability a classifier to not label a true negative observation as positive.

True Positive

True Positive + False Positive

ChrisAlbon

PRINCIPAL COMPONENTS

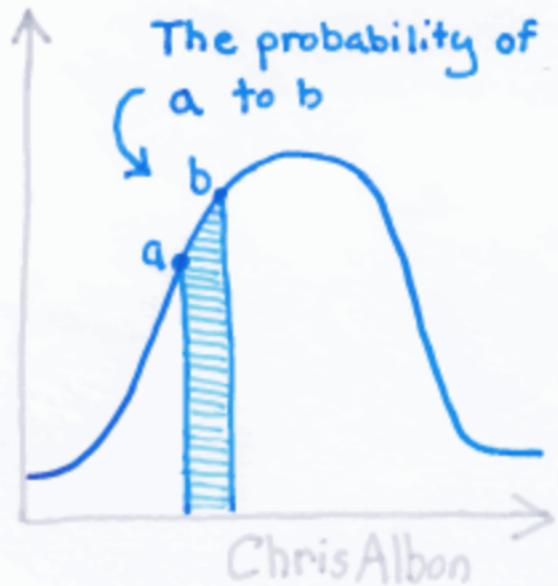
Principal components are the linear combination of features that have the maximum variance out of all linear combinations.

Alternative interpretation: Principal components are low dimensional linear surfaces closest to the observations.

PDF

Probability Density Function

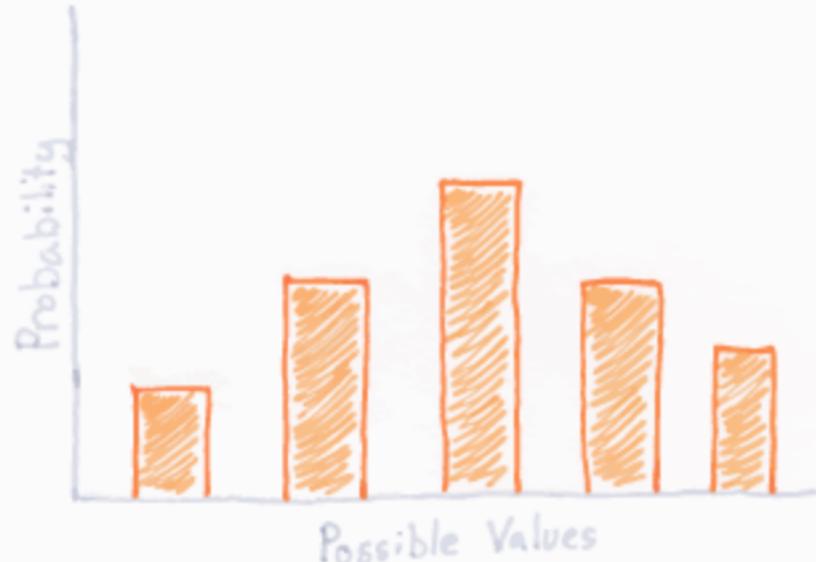
The PDF is the probability distribution of a continuous random variable. PDFs tell us the probability of an infinitely small region. We can use integration to find the probability.



PMF

Probability Mass Function

The probability distribution of a discrete random variable.



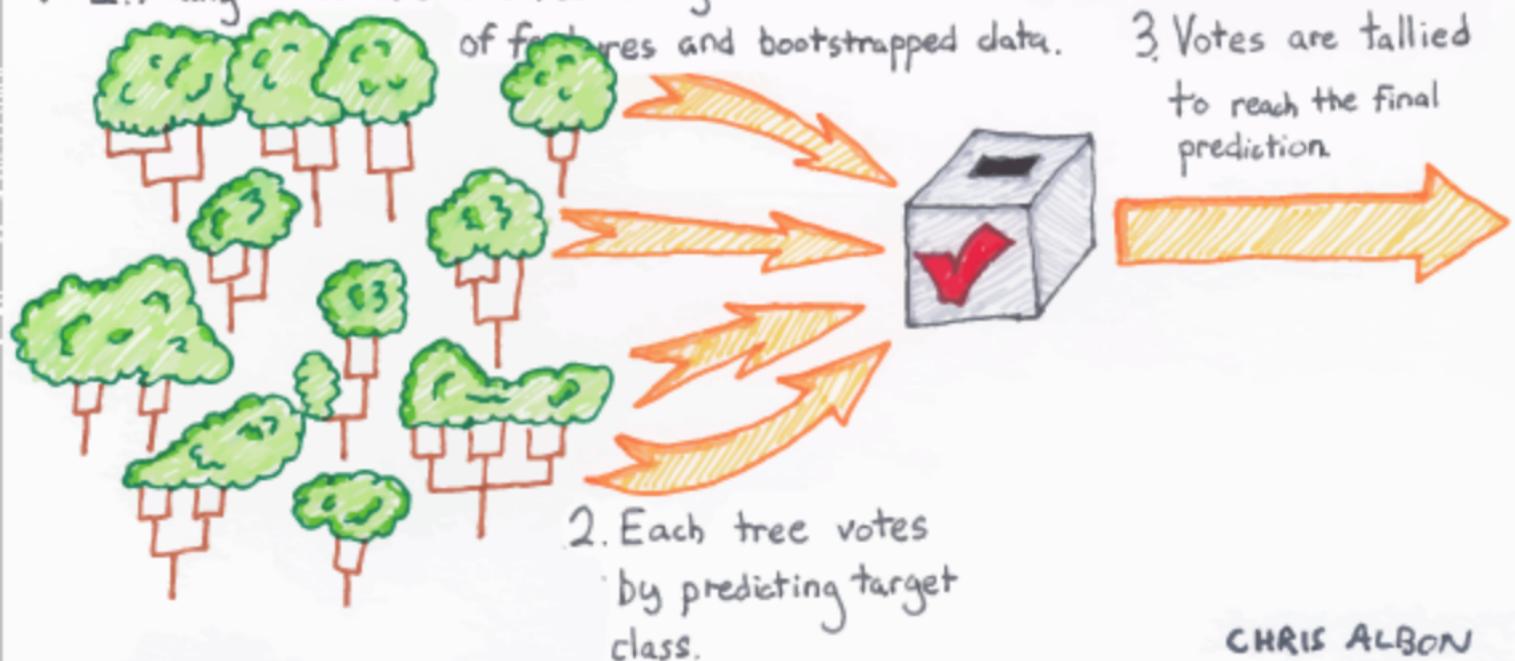
ChrisAlbon

RANDOM FOREST

1) Many trees are created using random subsets of features and bootstrapped data.

CLASSIFICATION

3. Votes are tallied to reach the final prediction.



RANDOM VARIABLE

Any variable taking on values randomly.

The values a random variable can take on is determined by a probability distribution

For example, the outcome of a dice roll is a random variable.

Chris Albon

RECALL

"Recall is about the real positives"

True Positives

True Positives + False Negatives

Recall is the ability of the classifier to find positive examples. If we wanted to be certain to find all positive examples, we could maximize recall.

Chris Albon

REGRESSION

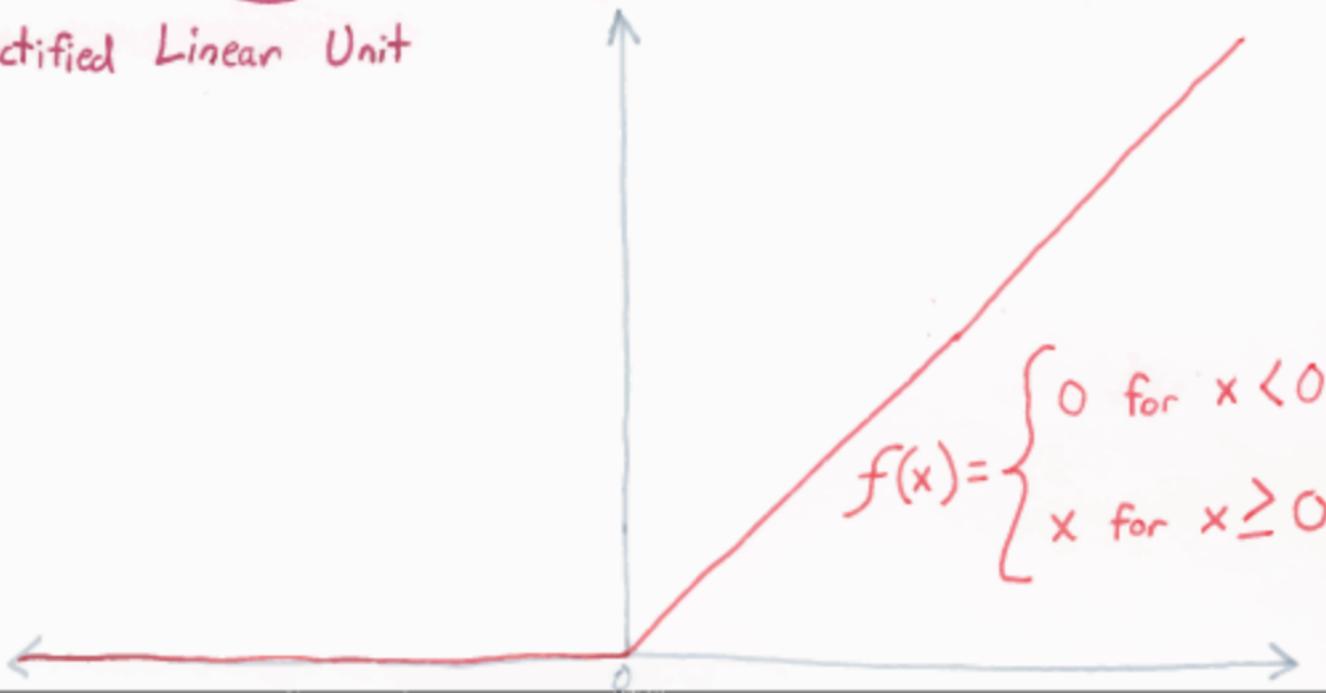
Regression trains models to predict quantitative targets. Example home price.



ReLU ACTIVATION FUNCTION

Rectified Linear Unit

BY CHRIS ALBON



RIDGE REGRESSION

Residual sum of
squares →

$$RSS + \lambda \sum_{j=1}^p \hat{B}_j^2$$

↑
Tuning parameter

Parameters squared

Shrinkage

Remember:
Standardize
the data first.

Chris Albon

Disadvantage:
parameters cannot
be zero like
with Lasso
regression.

SENSITIVITY

Also called "recall"

Everything predicted
positive correctly

True Positives

True Positives + False Negatives

Everything
actually positive

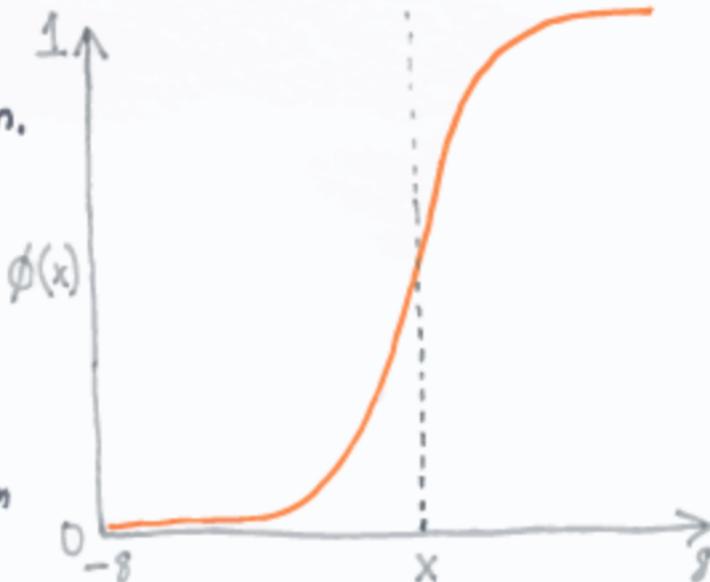
SIGMOID

ACTIVATION FUNCTION

Often used in the output unit for binary classification.

$$\phi(x) = \frac{1}{1+e^{-x}}$$

Returns a value bounded between 0 and 1.



SOFTMAX

NORMALIZATION

Reduces the influence of outliers without having to drop them.

$$x_i' = \frac{1}{1 + e^{-(x_i - \bar{x})/\sigma}}$$

Normalized value Euler's Number mean Standard deviation

SOURCES OF UNCERTAINTY

1. Inherent randomness in the universe.

Example: Quantum mechanics.

2. Inability to completely observe a phenomena, even when it is deterministic.

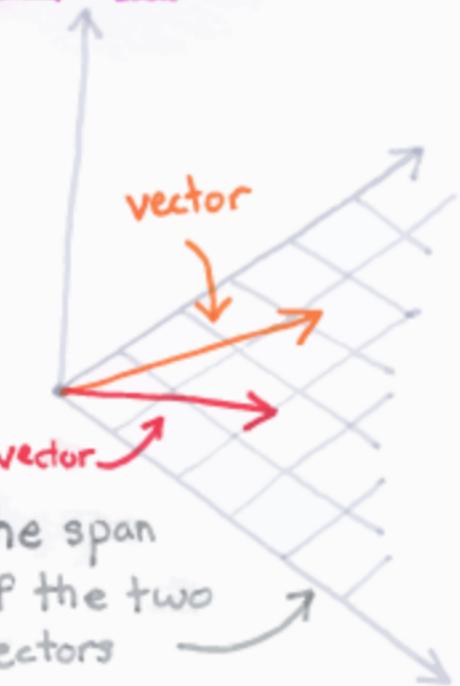
Example: Observing crime.

3. Inability to perfectly model a phenomena.

Example: Models predicting crime are simplifications.

SPAN

The span of a matrix is the set of all points reachable using a linear combination of the vectors of that matrix.



STANDARD ERROR OF THE MEAN

Standard error of the mean is the standard deviation of the sample mean. Estimated by dividing the Sample Standard deviation by Square root of the number of observations.

Sample standard deviation $s \leftarrow$

$$SEM = \frac{s}{\sqrt{n}}$$

Number of observations

STEMMING WORDS

Stemming reduces a word to its stem. The result is less readable by humans but makes the text more comparable across observations.

EXAMPLE: "Tradition" and "Traditional" have the same stem: "tradit"

ChrisAlbon

STOP WORDS

Any word to remove before processing. Frequently Stop words are extremely common Words with little informational value.

Examples

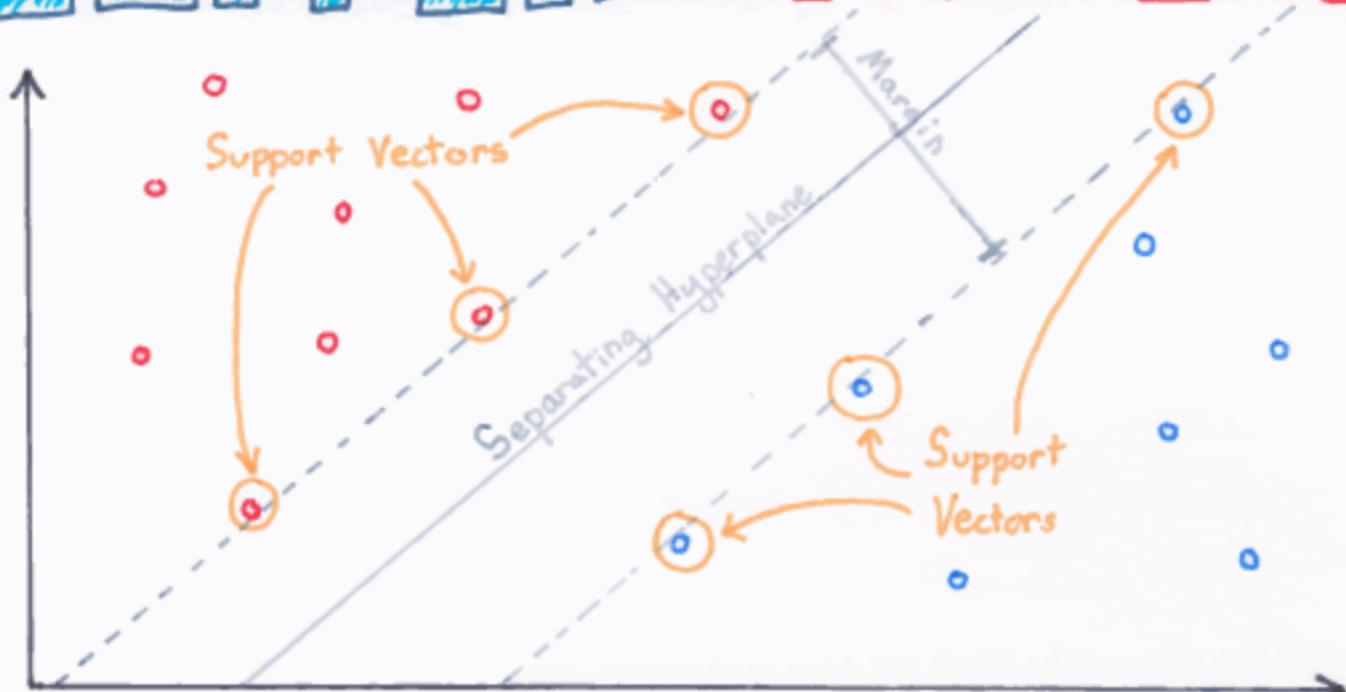
- it
- me
- myself
- we
- the
- and

STRATEGIES FOR HIGHLY IMBALANCED CLASSES

1. Collect more data
2. Choose a loss function suited for imbalanced classes like precision or recall.
3. Weight the classes
4. Downsampling and upsampling.

ChrisAlbon

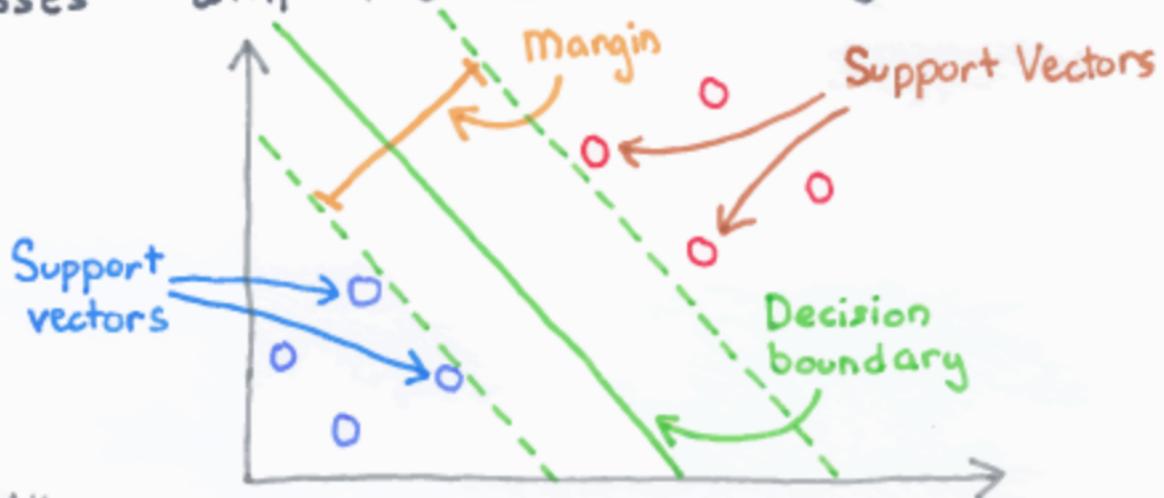
SUPPORT VECTORS



ChrisAlbon

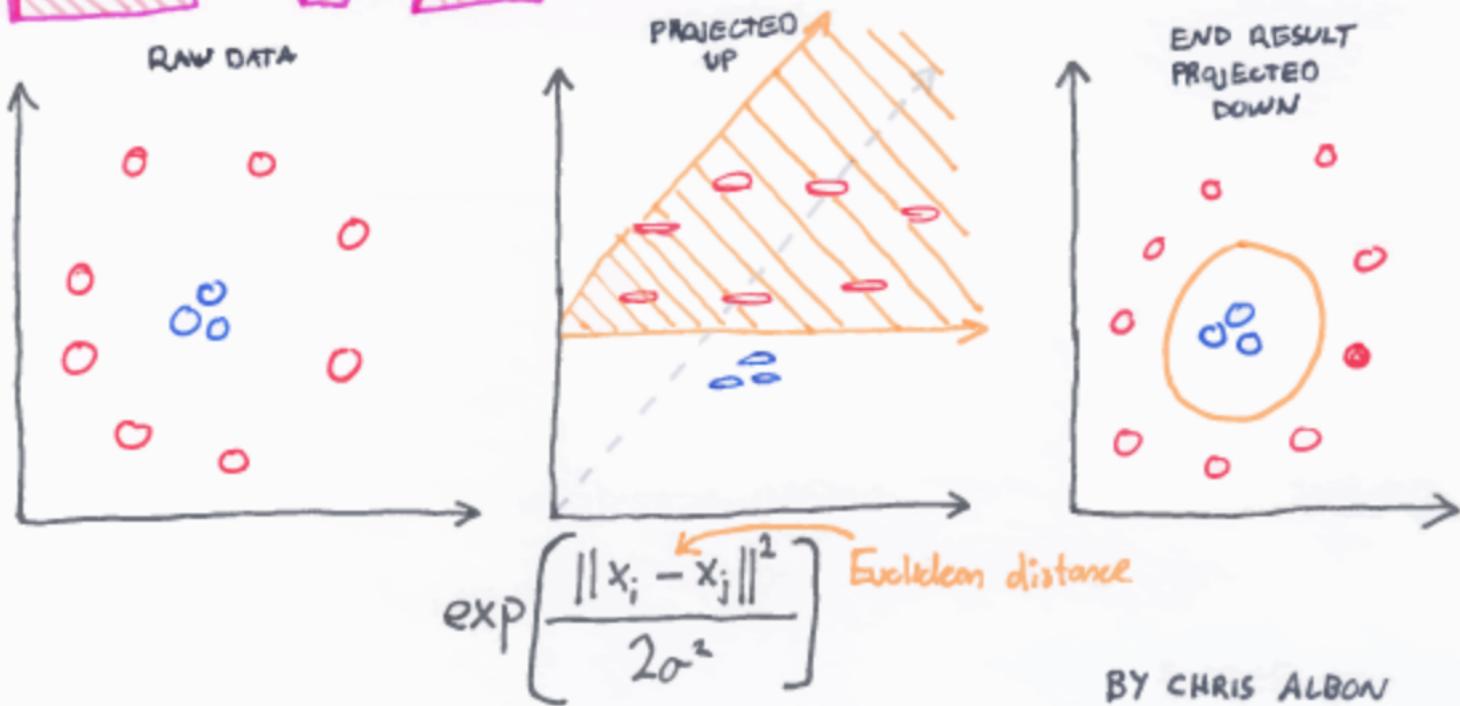
SVC

Finds the linear hyperplane that separates classes with the Maximum Margin.



EVF

RADIAL BASIS FUNCTION KERNEL



TENSORS

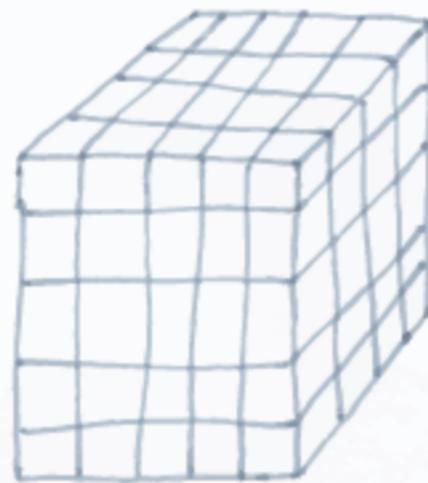
An n-dimensional array of numbers
arranged in a "grid".



1-dimension
tensor.
(vector)



2-dimension
tensor.
(matrix)



3-dimension
tensor.

TF-IDF

TF-IDF is a measure of originality of a word by comparing the number of times a word appears in a doc with the number of docs the word appears in.

$$\text{TF-IDF} = \text{TF}(t, d) \times \text{IDF}(t)$$

↑
Term frequency
↑
Number of times term t
appears in a doc, d

↑
Inverse document
frequency

$$\log \frac{1 + \frac{n}{\text{# of documents}}}{1 + \text{df}(d, t)} + 1$$

↑
Document frequency
of the term t

THEREFORE :
BECAUSE

NOTATION

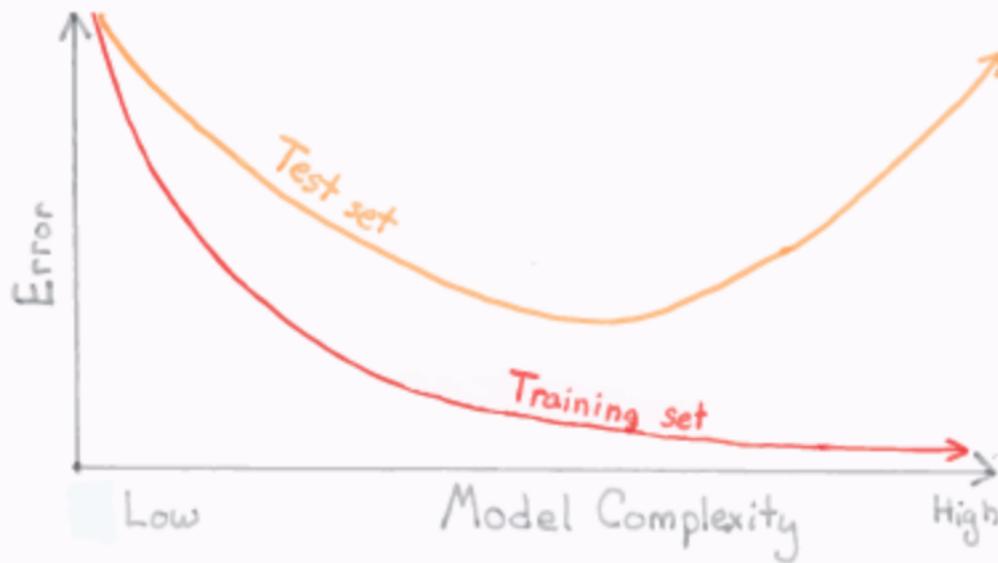
+ THEREFORE
BECAUSE

THE EFFECT OF **DROPOUT** ON HIDDEN UNITS

In dropout, each hidden unit must learn to perform well regardless of the other units in the network. The learned robustness helps the network perform well in the face of unseen test data.

Chris Albon

THE EFFECT OF MODEL COMPLEXITY TRAINING AND TEST ERROR



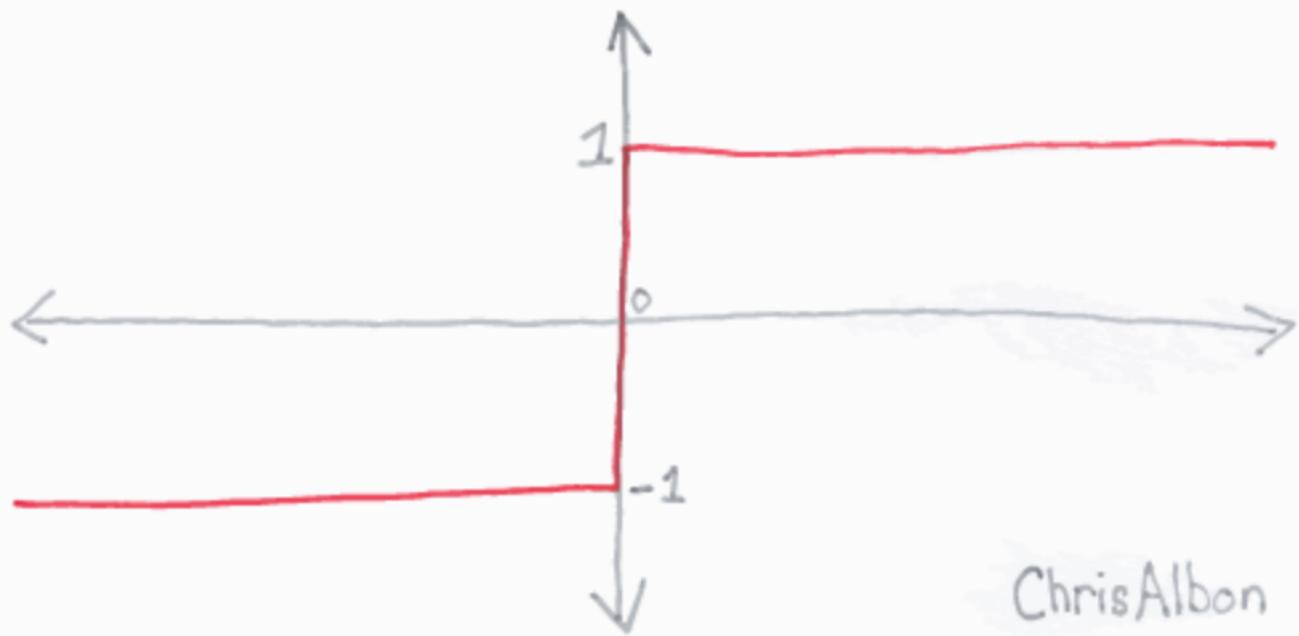
ChrisAlbon

THE RANDOM IN RANDOM FOREST

1. Each tree gets a random sample of observations with replacement.
2. Each tree gets all features, but at each node only a subset of those features are available.

Chris Albon

THRESHOLD ACTIVATION



ChrisAlbon

TOKENIZING TEXT

Splitting up text into individual units like paragraphs, sentences, or words.

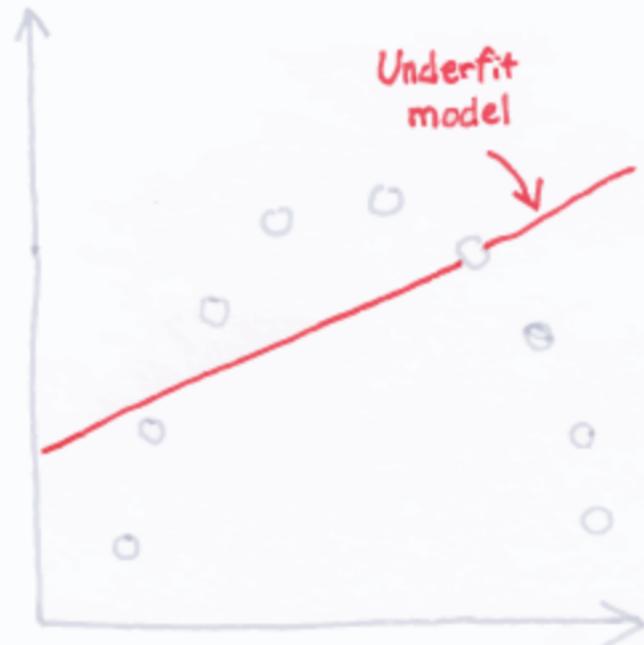
Example:

"I like birds" → "I", "like", "birds"

ChrisAlbon

UNDERFITTING

A model is underfit when it fails to capture the pattern in the data. It suffers from high bias.



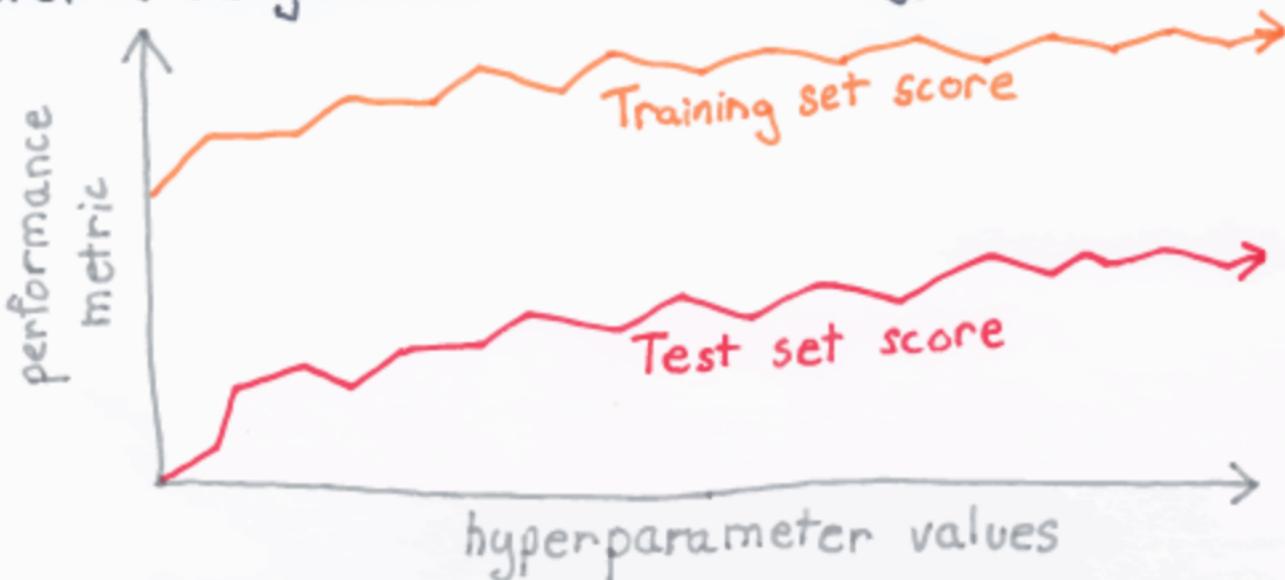
UNDERFLOW

Underflow occurs when a number is so small that it is too small to be represented by the computer. The Computer will most often round these values to zero, which can be problematic because zero often behaves differently to small numbers.

Chris Albon

VALIDATION CURVE

Validation curves visualize the performance metric over a range of values for some hyperparameter.



VARIANCE

$$\text{Var}[\hat{f}(x)] = E[\hat{f}(x)^2] - E[\hat{f}(x)]^2$$

Variance is the amount our predicted values would change if we had a different training dataset. It is the "flexibility" of our model, balanced against bias.

VIF

VARIANCE INFLATION FACTOR

Measures the effect of collinearity among features.

Specifically, measures how much the variance of a model parameter increases if features are correlated.

To calculate VIF we make the feature

the target of the model. Then run

the model and calculate the R^2 :

$VIF = 1 / (1 - R^2)$
 $VIF = 1$, not correlated. $VIF > 3$ correlated.

$$VIF_i = \frac{1}{1 - R_i^2}$$

ChrisAlbon

VECTORES

Vectors are ordered arrays of numbers (scalars).
Can be thought of as both data and geometrically:

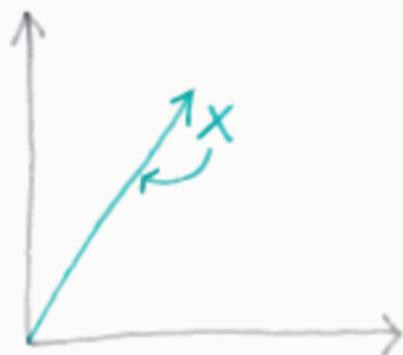
scalar
DATA:

$$x = [x_1 \ x_2 \dots \ x_n]$$

number of elements

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$

GEOMETRICALLY:



ZERO-ONE LOSS

Indicator Function

$$L_{0-1}(y_i, \hat{y}_i) = I(\hat{y}_i \neq y_i)$$

True class Predicted class True class