



Accuracy

$$\text{Accuracy} = \frac{1}{n} \sum \mathbf{1}(\hat{y}_i = y_i)$$

↳ number of observations

Predicted y ↳ True y
↓ ↗
 $\mathbf{1}$ (Indicator function)

Common metric in classification. Fails when in the presence of highly imbalanced classes. In those case F1 score is likely more appropriate.

BY CHRIS ALBON

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.

ADJUSTED R²

Intuition: Once all the correct features have been added, additional features should be penalized.

$$\hat{R}^2 = 1 - \frac{\text{Residual Sum of Squares}}{\text{Total Sum of Squares}} / \left(\frac{\text{Number of Features}}{\text{Number of Observations}} - 1 \right)$$

$\frac{\text{RSS}/(n-d-1)}{\text{TSS}/(n-1)}$

ChrisAlbon

ACCOLOMERATIVE CLUSTERING

All observations start as their own cluster. Clusters meeting some criteria are merged. This process is repeated, growing clusters until some end point is reached.

Chris Albon

AIC

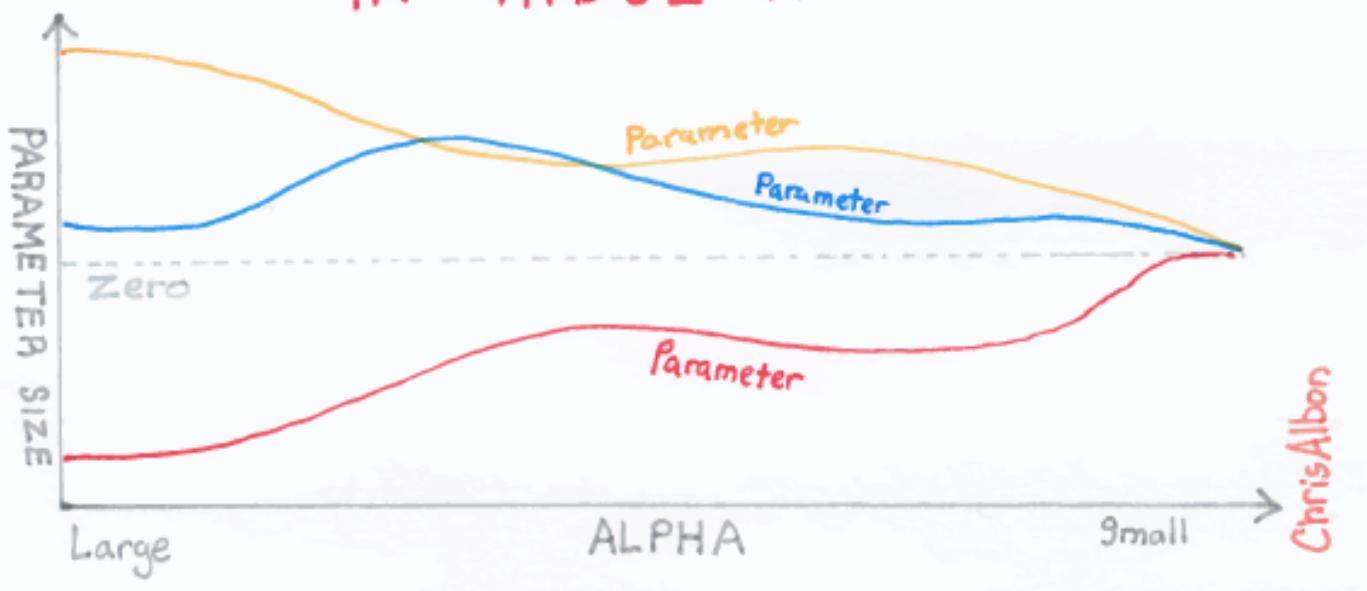
"Ey-Kye-Ih-Key"
Information Criteria

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

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

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

ALPHA IN RIDGE REGRESSION



AVOID OVER-FITTING

Simple Models

Cross-validated Evaluation

Regularization

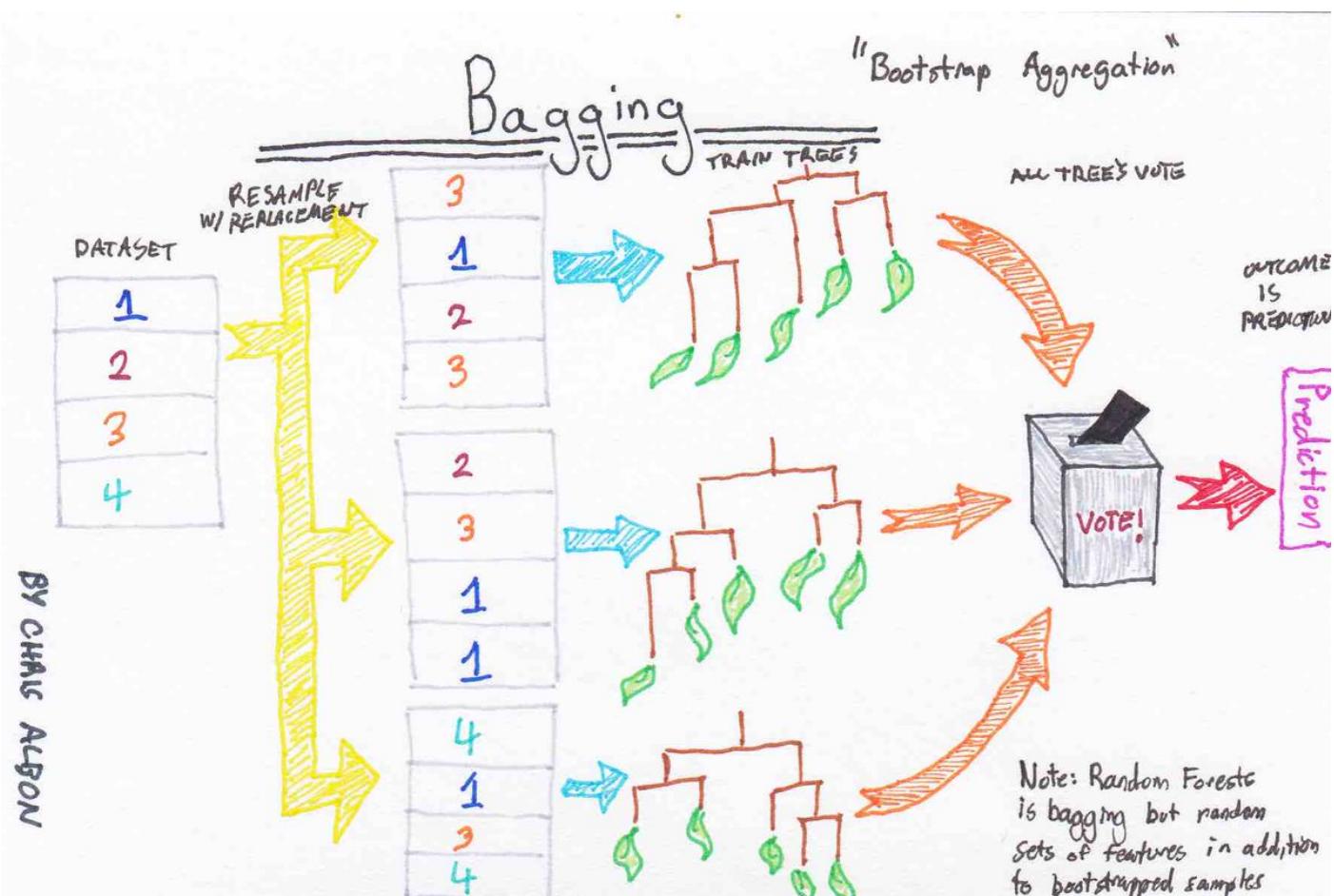
Get More Data

Ensemble Learning

BACKWARD STEPWISE SELECTION

1. Create a model, m_0 , with all features, f .
2. For $K=F \ K-1, \dots, 1$:
 - a) Create a model with one less feature from the last model.
 - b) Repeat for all features
 - c) Choose model with best evaluation metric and define as m_{K-1}
3. Select the best model from m_0 to m_f using CV.

BY CHRIS ALBON



DEEPLERNING

BASIC PARTS

1. Data.
2. Loss function. Example: cross-entropy.
3. Optimization algorithm. Example: Adam
4. Network architecture. Example: Dense layers
5. Test data.
6. Evaluation metric. Example: Accuracy

Chris Albon

Bayes THEOREM

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

Bias - Variance Tradeoff

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

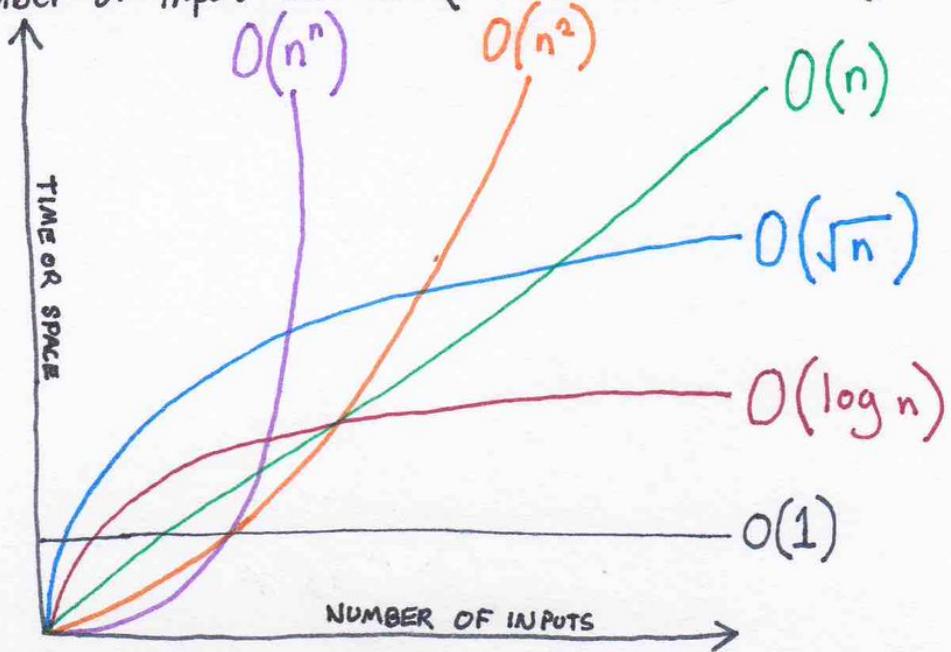
↓ ↓ ↓ ↓
predicted true predicted average predicted value irreducible error

Bias²
How much predicted values differ from true values.

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

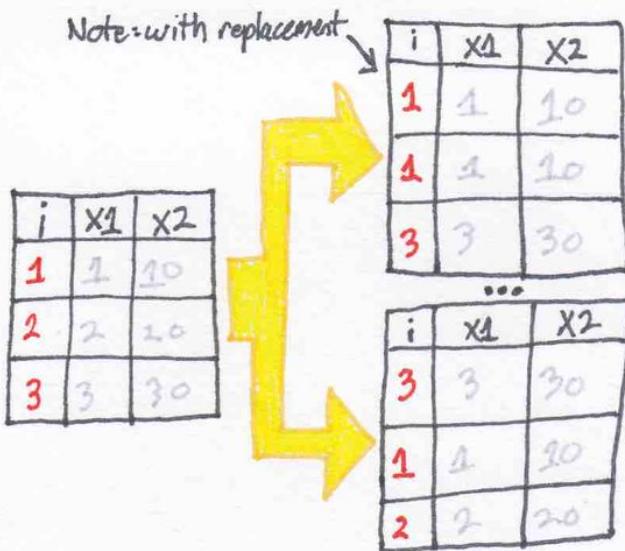
BIG O

Representation of how an algorithm's space requirement or run time increases as the number of inputs increases (in worst case scenario).



BOOTSTRAP

Note: with replacement



Bootstrap allows us to simulate obtaining many new datasets by repeated sampling with replacement from the original original dataset.

Note: normally we would create many more than two bootstrapped data sets

Brier Score

$$BS = \frac{1}{n} \sum_{i=1}^n (P_i - O_i)^2$$

↑
predicted probability
↓
number of observations

↑
actual outcome
↓

Brier score shows the squared mean difference between the predicted probability of all observations with their actual outcome.

The lower the Brier Score, the better. Ranges between 0 and 1.

BY CHRIS ALBON

C

INVERSE OF REGULARIZATION STRENGTH

A hyperparameter, α , is used to control the regularization strength. For example:

$$\alpha \sum_{j=1}^n |\hat{\beta}_j|$$

L1 norm

However, often we will see the hyperparameter C which is the inverse

$$C = \frac{1}{\alpha}$$

VIF

- VIF measures the effect of collinearity ^{among} ~~among~~ features.
- VIF measures how much the variance of a model parameter (coefficient) increases if features are correlated.
- To calculate VIF we make a feature the target of the model (instead of the dependent variable). Run the model, then calculate the R^2 :

$VIF_i \leq 1 \rightarrow$ none to assess

$VIF_i \leq 5 \rightarrow$ moderate

$VIF_i \leq 10 \rightarrow$ high

$$VIF_i = \frac{1}{1 - R_i^2} \quad \text{where } i \text{ is the } i^{\text{th}} \text{ predictor.}$$

BY CHRIS ALBON

IF TWO VECTORS ARE
INDEPENDENT, PROB.
~~BEAN~~ CATEGORY OF

Chi-Squared For Feature Selection

FEATURE IN TARGET

CLASS A WILL BE THE

SAME AS TARGET CLASS B.

$$\chi^2 = \sum_{i=1}^n \frac{(O_i - E_i)^2}{E_i}$$

We calculate between
each feature and target
vector. If features are
independent, then feature
is likely irrelevant.

count of observations in ~~target~~
group/class i

O_i

E_i

expected number of observation
in ~~target~~ class i

Chi-SQUARED

$$\chi^2 = \sum_{i=1}^n \left(\frac{O_i - E_i}{E_i} \right)^2$$

O_i = observed value

E_i = expected value

Tests if categorical data shows up at a rate different than random.

BY CHRIS ALBON

Coefficient of Determination, R^2

$$R^2 = 1 - \frac{\sum_i (y_i - \hat{y}_i)^2}{\sum_i (y_i - \bar{y})^2}$$

Annotations:

- True y (green arrow pointing to y_i)
- Predicted y (green arrow pointing to \hat{y}_i)
- Mean y (red arrow pointing to \bar{y})
- How much of the variance in the target vector is explained by the features (red text)

Common Mathematical Notation

\neg	NOT	\perp	orthogonal, independent
$a \in b$	a is an element of b	\propto	proportional
\exists	such that	\sum	summation
\therefore	therefore	\prod	product
\because	because	\lfloor	floor (round down to nearest integer)
\Rightarrow	then	\lceil	ceiling (round up to nearest integer)
\Leftrightarrow	if and only if	$x y$	x , given y
\exists	there exists	$\max()$	max value for set/list
\forall	for all	$\underset{x}{\operatorname{argmin}} f(x)$	value of x the minimizes $f(x)$
\parallel	parallel		

COMMON OPTIMIZERS FOR NEURAL NETS

1. Stochastic gradient descent
2. Stochastic gradient descent with momentum.
3. RMSProp
4. Adam

Chris Albon

COMMON OUTPUT LAYER

ACTIVATION FUNCTIONS

BINARY CLASSIFICATION: Sigmoid

MULTICLASS CLASSIFICATION: Softmax

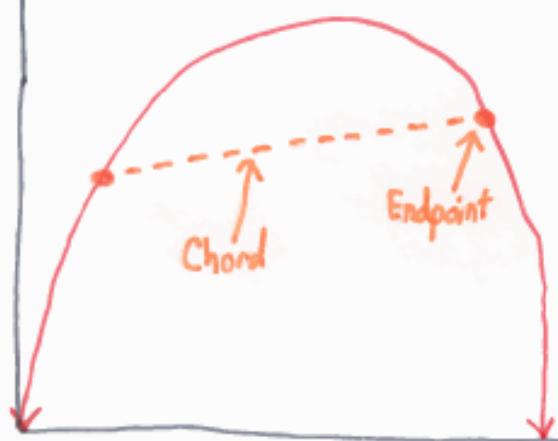
REGRESSION: No activation function

Chris Albon

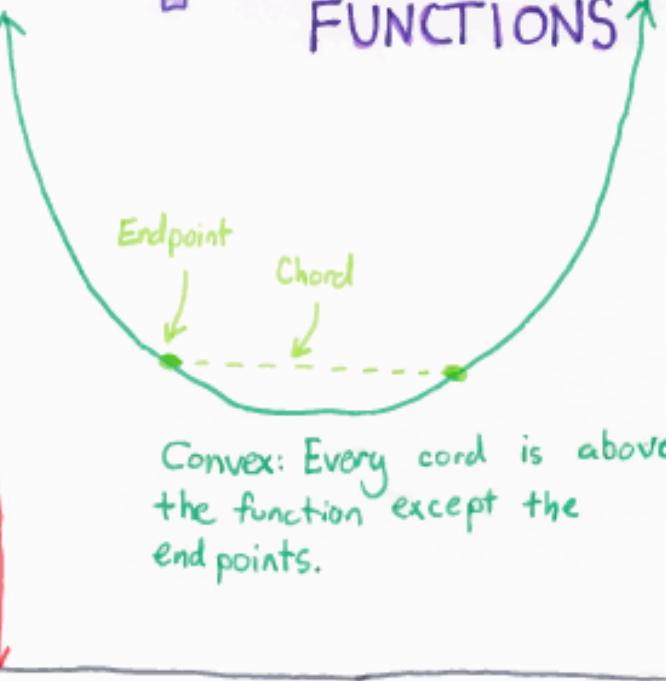
CONCAVE & CONVEX

FUNCTIONS

↑ Concave: Every chord is below the function except the endpoints.



↓ Convex: Every chord is above the function except the end points.



BY CHRIS ALBON

Conditional Probability

$$P(A \text{ and } B) = P(A) \times P(B | A)$$

EXAMPLE:

$$P(A \text{ and } B) = P\left(\frac{\text{DRAW ACE}}{52} \text{ and } \frac{\text{DRAW KING}}{51}\right) = \left(\frac{4}{52}\right) \times \left(\frac{4}{51}\right)$$

BY CHRIS ALBON

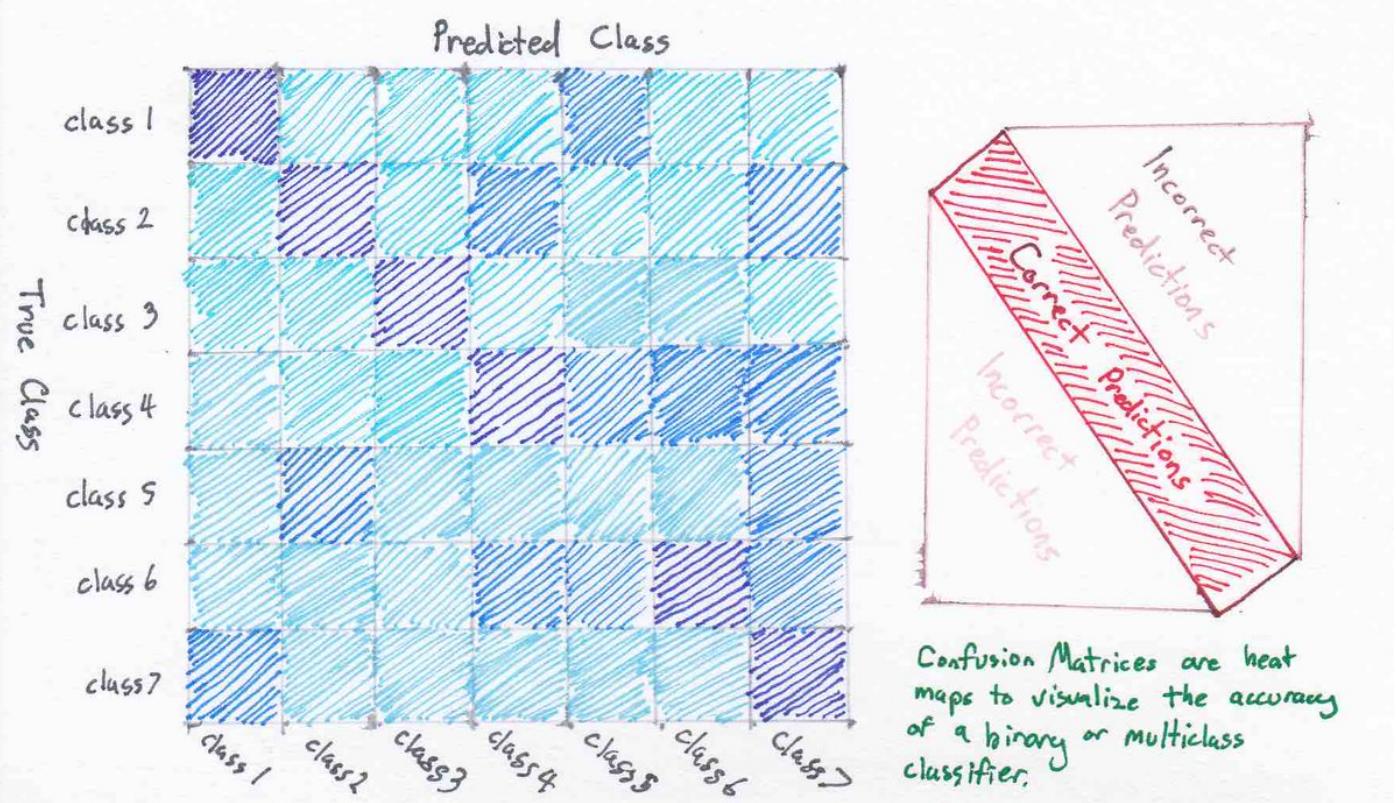
CONDITIONING

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.

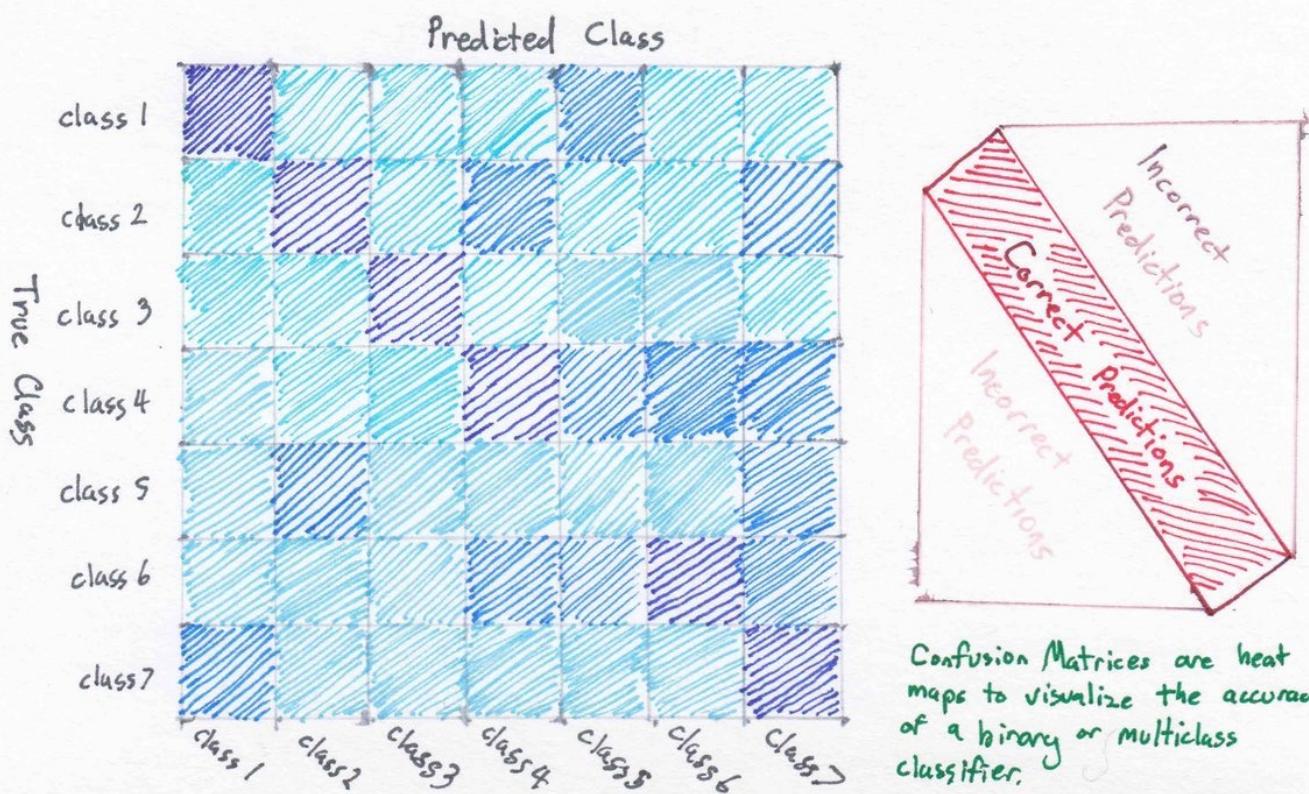
Chris Albon

Confusion Matrix

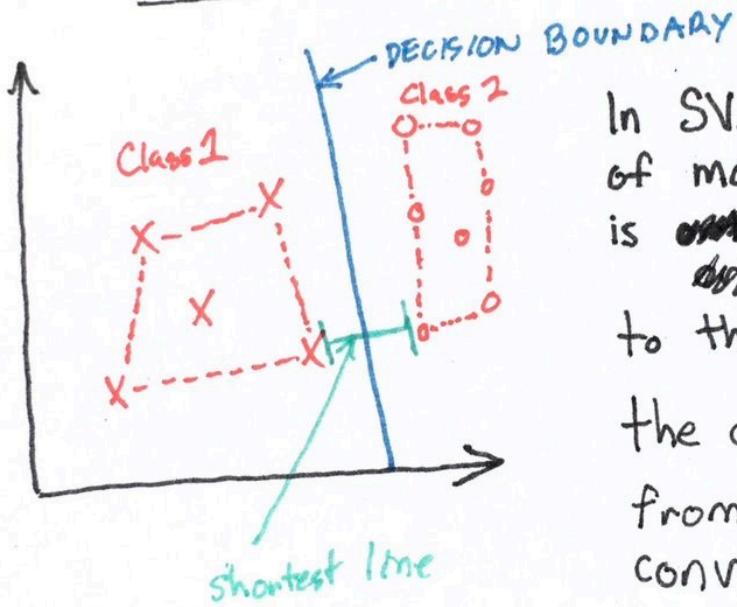
BY CHRIS ALBON



Confusion Matrix



Convex Hull's Relation To Support Vector Machines



In SVM the line of maximum margin is orthogonal to the line connecting the closest two points from each group's convex hull.

Correlation of Error Terms

If there is correlation in the error terms,
we will have underestimated STANDARD ERRORS.

HETERO~~SCE~~DASTICITY

Breaks Gauss - Markov Theorem:

NOT BLUE. CAUSES ESTIMATES OF OLS
VARIANCE TO BE BIASED.

Correlation

I'm drunk

$$\frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

COST LOSS ERROR

MINIMIZATION

Cost functions, loss functions, and error functions mean the same thing. They are the objective function we are trying to train a model that minimizes or maximizes, depending.

BY CHRIS ALBON



Cost functions, loss functions, and error functions mean the same thing. They are the objective function we are trying to train a model that minimizes or maximizes, depending.

Covariance and Correlation

Correlation

$$\left\{ \begin{array}{l} \frac{E(x_i - \bar{x})(y_i - \bar{y})}{(\sigma_x \sigma_y)} \leftarrow \text{covariance} \\ \end{array} \right.$$

$$\underline{C_p}$$

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

↓ ↓
residual sum estimated error variance
of squares number of features

$2d\hat{\sigma}^2$ is a penalty to adjust for
the fact that the training data
Underestimates test error.

C_p is used in model selection to
Compare different models.

CROSS-ENTROPY

Proportion of observations in mth region of k class.

$$D = - \sum_{k=1}^K \hat{P}_{mk} \log \hat{P}_{mk}$$

Region Class Region Class

The closer D is to 0 the purer the classes.

$0 \leq -\hat{P}_{mk} \log \hat{P}_{mk}$, so the larger the value the less pure.

BY CHRIS ALBON

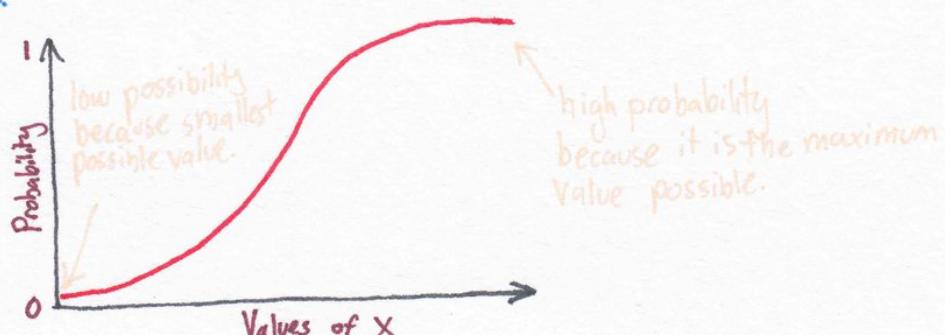
Cumulative Distribution Function

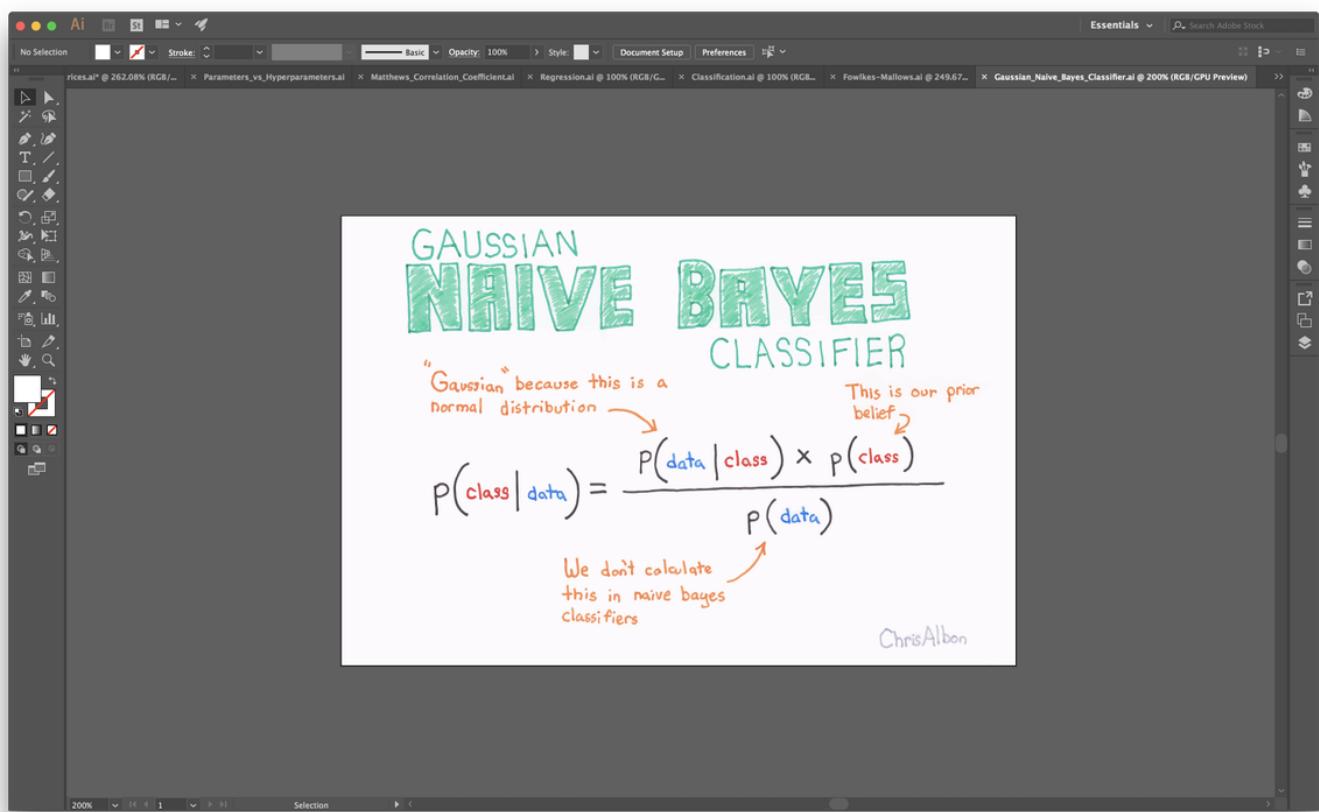
$$f(x) = \Pr[X \leq x]$$

where $x \in \mathbb{R}$

CDF tells us the probability a random variable/function returns a value less than some specified value.

CDF is the accumulation of the ~~possible~~ probability of values up to some value.





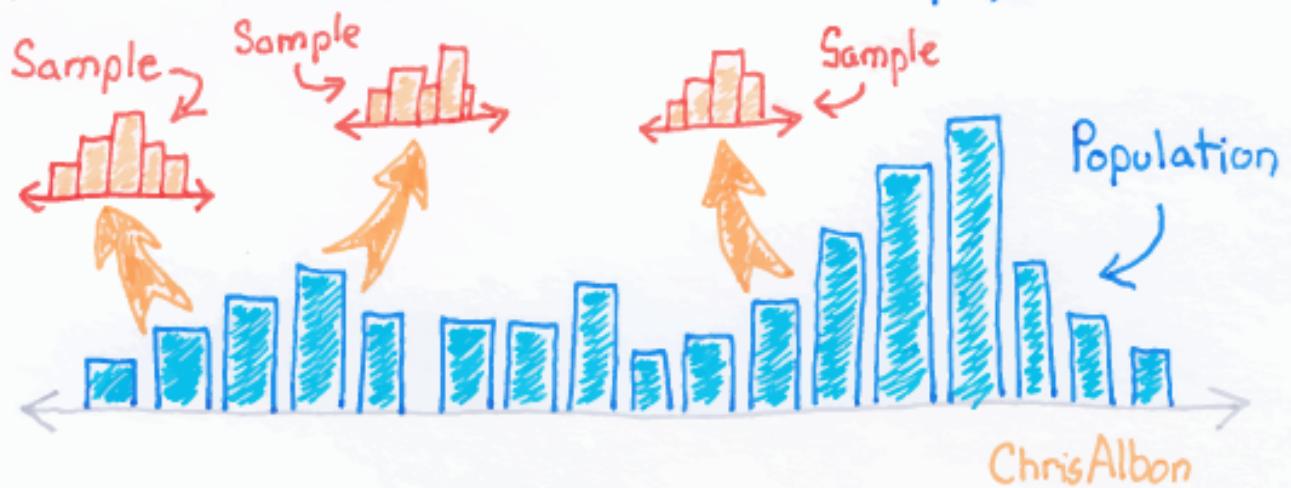
Curse of **DIMENSIONALITY**

As the dimensionality of the features space increases, the number of configurations can grow exponentially, and thus the number of configurations covered by an observation decreases.

Chris Albon

DATA-GENERATING DISTRIBUTION

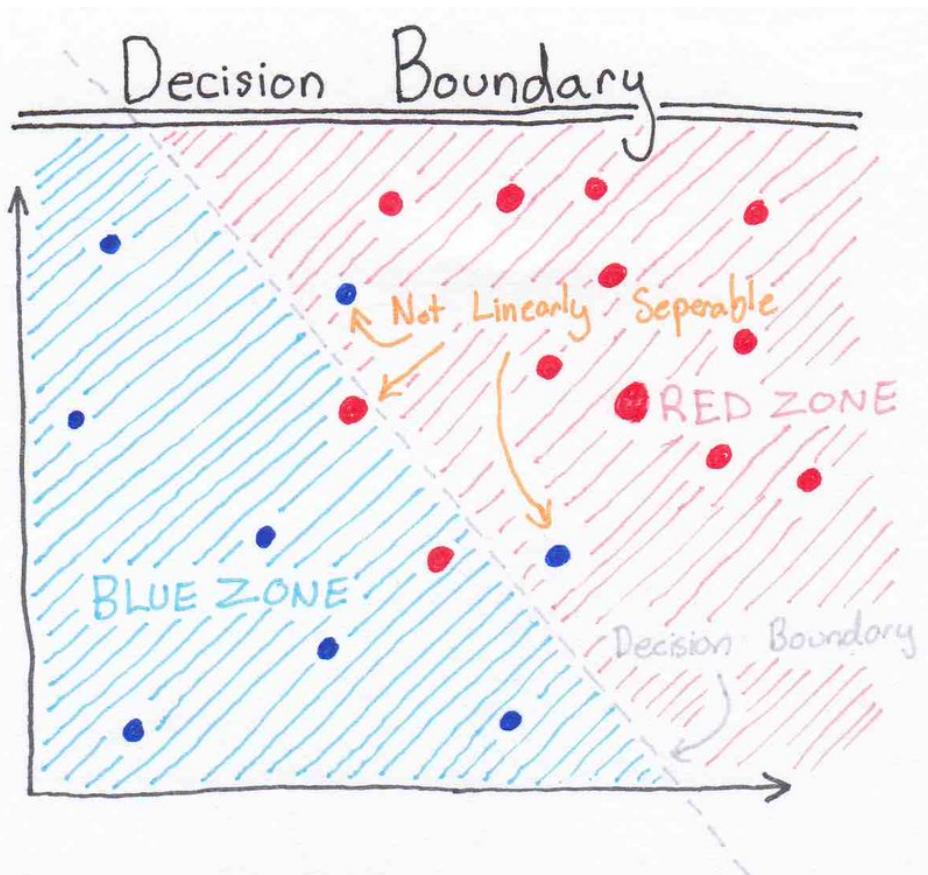
The unseen distribution of the population.



DATASET AUGMENTATION

- Often used in image recognition.
- Creates noise by rotating, scaling, shifting etc...
images in computervision problems. Other
ways of injecting noise are used in other problems.
- Can greatly reduce generalization error.

Chris Albon



BY CHRIS ALBON

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

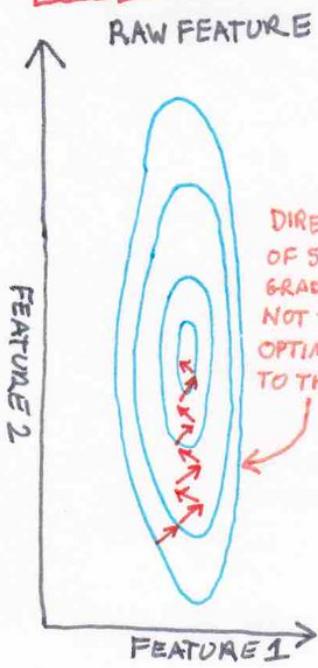


Downsampling

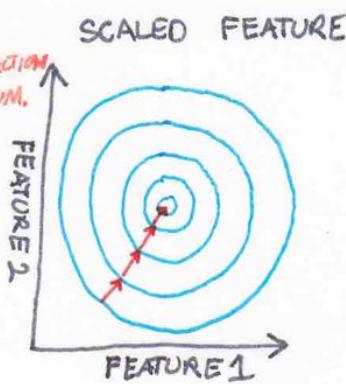
A strategy to handle imbalanced classes by creating a random subset of the majority of equal size to the minority class.

ChrisAlbon

EFFECT OF FEATURE SCALING ON GRADIENT DESCENT



Gradient descent will take longer if features are not similarly scaled. Andrew Ng's rule is that if x_i ranges from -3 to 3 or $-\frac{1}{3}$ to $\frac{1}{3}$ then it is fine.



Effect of K's value in K-nearest neighbors

Small K = Low Bias, High Variance

Large K = High Bias, Low Variance

EIGENVECTOR

$$A v = \lambda v$$

Diagram illustrating the relationship between a square matrix A , a vector v , and a scalar λ :

- A blue arrow labeled "Square matrix" points to the matrix A .
- A green arrow labeled "eigenvalue" points to the scalar λ .
- A red arrow labeled "eigenvector" points to the vector v .

Chris Albon

ELASTIC NET

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

$$\text{RSS} + \alpha \rho \left(\frac{\|w\|_1}{\text{Weights}} + \frac{\alpha(1-\rho)}{2} \|w\|_2^2 \right)$$

Alpha determines the regularization strength.

Rho determines the balance between L1 and L2.

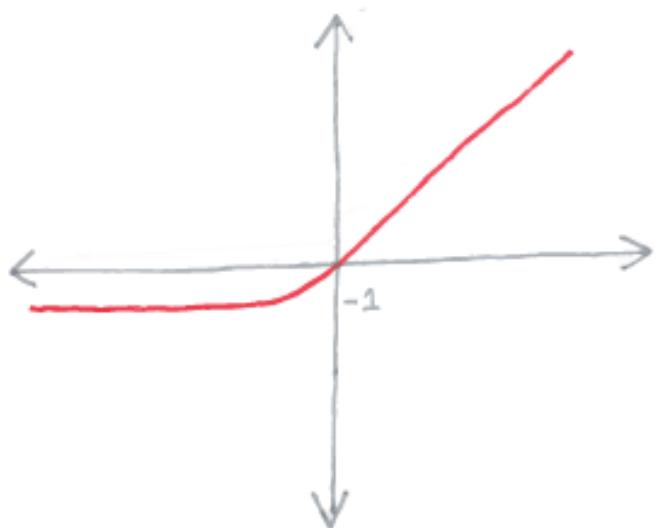
ChrisAlbon

ELUs

Exponential Linear Units

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

positive hyperparameter



Chris Albon

Error Types

TYPE
FALSE
Positive

TYPE
FALSE
NEGATIVE

Euclidean Distance

For vectors, p_i and q .

$$\text{Euclidean Distance} = \sqrt{\sum_{i=1}^n (q_i - p_i)^2}$$

Euclidean Norm

[Often used in rescaling via normalizing
observation values.]

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

Also called: L₂ norm

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

BY CHRIS ALBON

EXPLODING GRADIENT PROBLEM

When the gradients of the loss function with respect to the parameters in the early layers are very large. Causes unstable learning and poor performance. One potential fix is gradient clipping.

Chris Albon

F-STATISTIC

$$F = \frac{(TSS - RSS)/p}{RSS/(n-p-1)}$$

Total Sum Of Squares Residual Sum Of Squares
Number of Features
Number of Observations

Used by the F-test to test if groups of features are jointly statistically significant.

Chris Albon

F1 Score

$$F_1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

F1 score can be interpreted as the harmonic mean of precision and recall. Values range from 0(bad) to 1(good).

BY CHRIS ALBON

F1 Score

$$F_1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

F1 score can be interpreted as the harmonic mean of precision and recall. Values range from 0 (bad) to 1 (good).

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.
4. Random Forest feature importance
5. Lasso

LINEAR REGRESSION

FINDING PARAMETERS

$$w = \left(X^T X \right)^{-1} X^T y$$

Parameter vector
Feature matrix
Transposed feature matrix
Target vector

Analytical solutions are most common. However, the cost of computing pseudo-inverse of $X^T X$ means that gradient descent can be better in large data.

Chris Albon

FOWLKES-MALLONS

SCORE

Evaluate clusters when ground-truth is available. Used to evaluate pairs of clusters. Higher score is better. Ranges between 0 and 1

$$FMI = \frac{TP}{\sqrt{(TP + FP)(TP + FN)}}$$

TP
↓
Fowlkes-Mallows
Index

TP
↓
True Positive
FP
↓
False Positive

FN
↓
False Negative

TP: # of pairs of observations are part of the same cluster in both y and \hat{y} .

FP: # of pairs of observations that are part of the same cluster in y but part of different clusters in \hat{y} .

FN: # of pairs of observations that are not part of the same cluster in y but part of the same cluster in \hat{y} .

BY CHRIS ALBON

FOWLKES-MALLOWS

SCORE

Evaluate clusters when ground-truth is available. Used to evaluate pairs of clusters. Higher score is better. Ranges between 0 and 1

$$FMI = \frac{TP}{\sqrt{(TP + FP)(TP + FN)}}$$

↓
Fowlkes-Mallows
Index

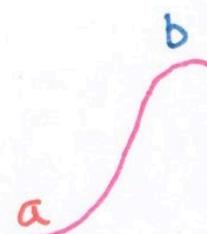
TP: # of pairs of observations are part of the same cluster in both y and \hat{y} .

FP: # of pairs of observations that are part of the same cluster in y but part of different clusters in \hat{y} .

FN: # of pairs of observations that are not part of the same cluster in y but part of the same cluster in \hat{y} .

True Positive
False Positive
False Negative

Fundamental Theorem Of Calculus


$$\text{area} = \boxed{F(b) - F(a)}$$

$= \frac{\text{anti-derivative}(b) - \text{anti-derivative}(a)}{F(b) - F(a)}$

General Additive Model

$$GAM = \hat{y}_i = B_0 + f(x_{i1}) + f(x_{i2}) \dots e_i.$$

- Every IV is a function.
- No interaction terms are possible.

GINI INDEX

Proportion of observations
in the m th leaf of K^{th} class.

$$G = \sum_{k=1}^K \hat{P}_{mk} \left(1 - \hat{P}_{mk} \right)$$

Diagram illustrating the components of the Gini Index formula:

- \hat{P}_{mk} : Proportion of observations in the m th leaf of K^{th} class. Arrows point from "leaf" to \hat{P}_{mk} and from "Class" to \hat{P}_{mk} .
- $1 - \hat{P}_{mk}$: Complement of the proportion above. Arrows point from "leaf" to $1 - \hat{P}_{mk}$ and from "Class" to $1 - \hat{P}_{mk}$.

Used at each node
to decide which
feature is best
to split on.

The smaller
the value of G
the more purity
there is in the
node.

Measure of purity
in tree based methods

GRADIENT DESCENT

Updated parameter

$$\tilde{w} \equiv w + \nabla w$$

Parameter Gradient Loss function

Negative $\rightarrow -\eta \nabla$ Learning rate

Parameter

Chris Albon

GREEK 1

LETTERS

A α alpha Δ δ delta

B β beta E ϵ epsilon

Γ γ gamma Z ζ zeta

Chris Albon

GREEK 3

LETTERS

N ν nu

Π π pi

Ξ ξ xi

Ρ ρ rho

Ο ο omicron

Σ σ sigma

HAMMING LOSS

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

Predicted class true class
 \hat{y}_i y_i
n number of observations

Chris Albon

HANDLING IMBALANCED CLASSES IN SUPPORT VECTOR MACHINES

In SVMs, the hyperparameter C determines the penalty for misclassification. To handle imbalanced classes we can weight C by class:

$$C_K = C * w_K \quad \begin{matrix} \text{penalty} \\ \text{for class } K \end{matrix} \quad \begin{matrix} \text{weight inversely proportional} \\ \text{to class } K \text{'s frequency} \end{matrix}$$

The goal is to prevent the majority class from overwhelming minority class.

Chris Albon

HANDLING OUTLIERS

1. If due to an error: drop, mark as missing value, mark as possible error.
2. If a legitimate but extreme value: decide if it is genuinely a member of the population we are trying to address with our model.

Chris Albon

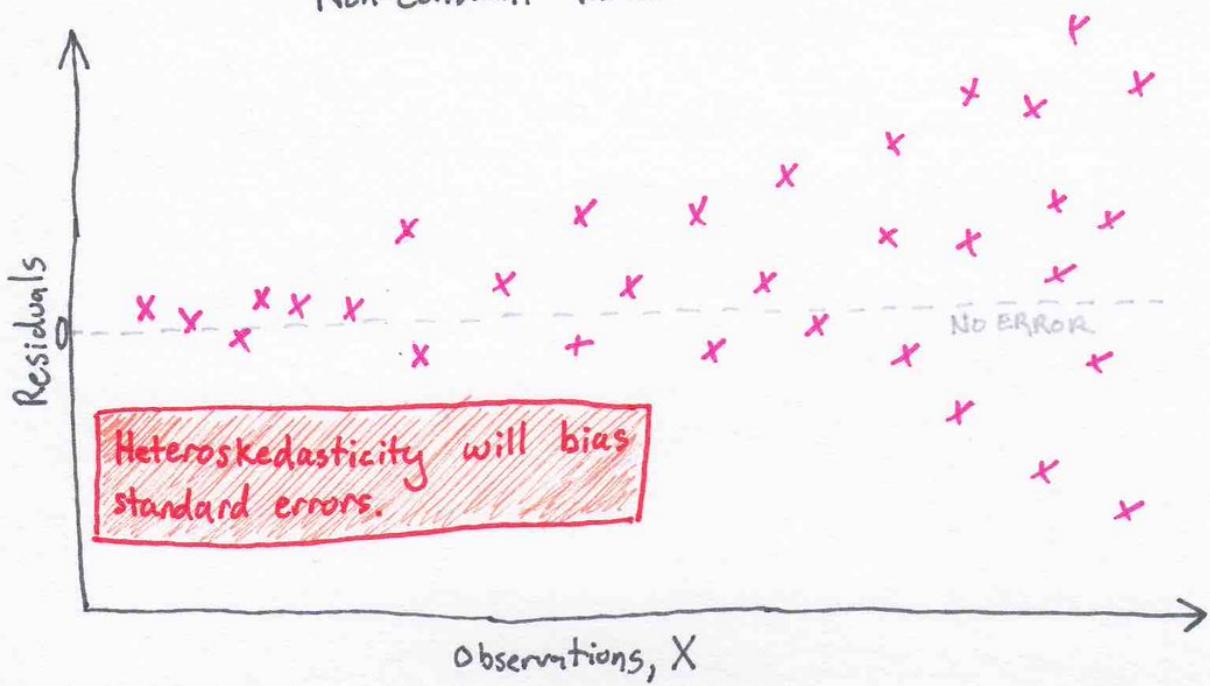
HESCIAN MATRIX

When we have a function taking in multiple inputs and returning one output, the square matrix of all its second order partial derivatives is called the Hessian matrix.

ChrisAlbon

HETEROSKEDASTICITY

Non-constant Variance in Errors.

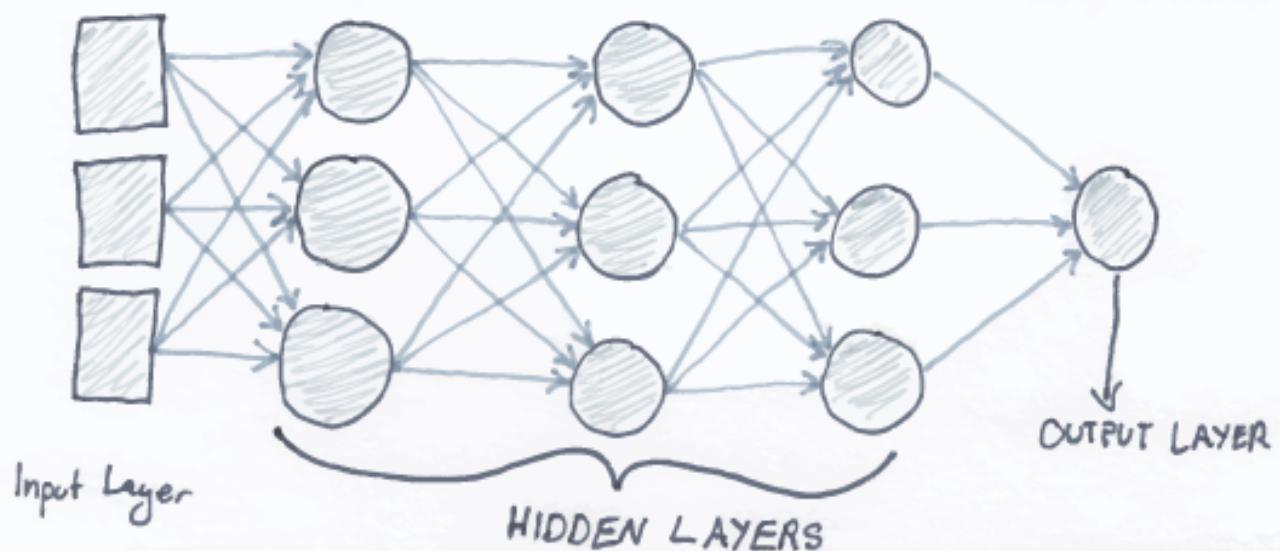


BY CHRIS ALBON

HIDDEN LAYER

VALUES

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



BY CHRIS ALDON

HINGE LOSS

Class labels are +1 and -1

$$L_{\text{Hinge}}(y, w) = \max \left\{ 1 - wy, 0 \right\}$$

Annotations:

- True y (green arrow pointing to y)
- Raw output of classifier (orange arrow pointing to wy)
- True y (green arrow pointing to the zero)

Used in support vector classifiers. When y and w have the same sign and $|w| \geq 1$ then Hinge Loss is 0. Otherwise, Hinge Loss increases as w increases.

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

How to avoid overfitting

Simple models

Cross-validation

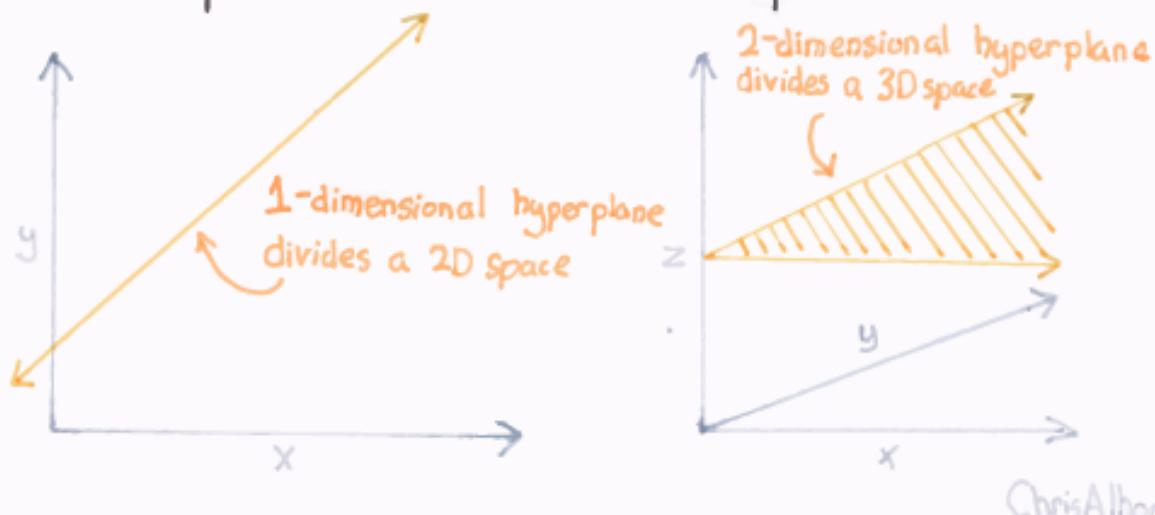
Regularization

Get More Data

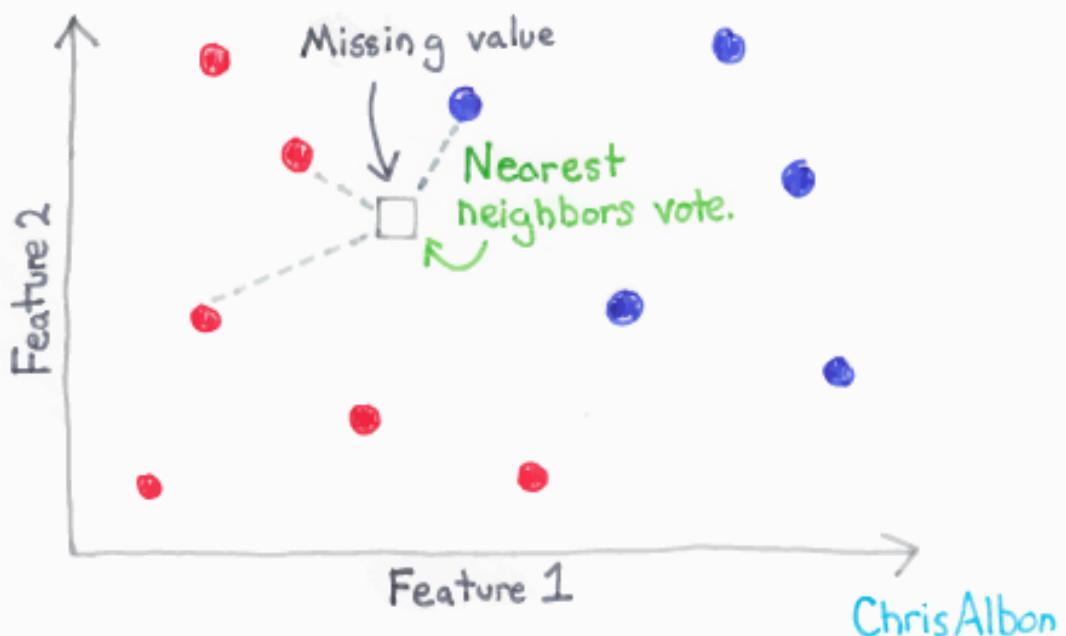
Ensemble Models

HYPERPLANE

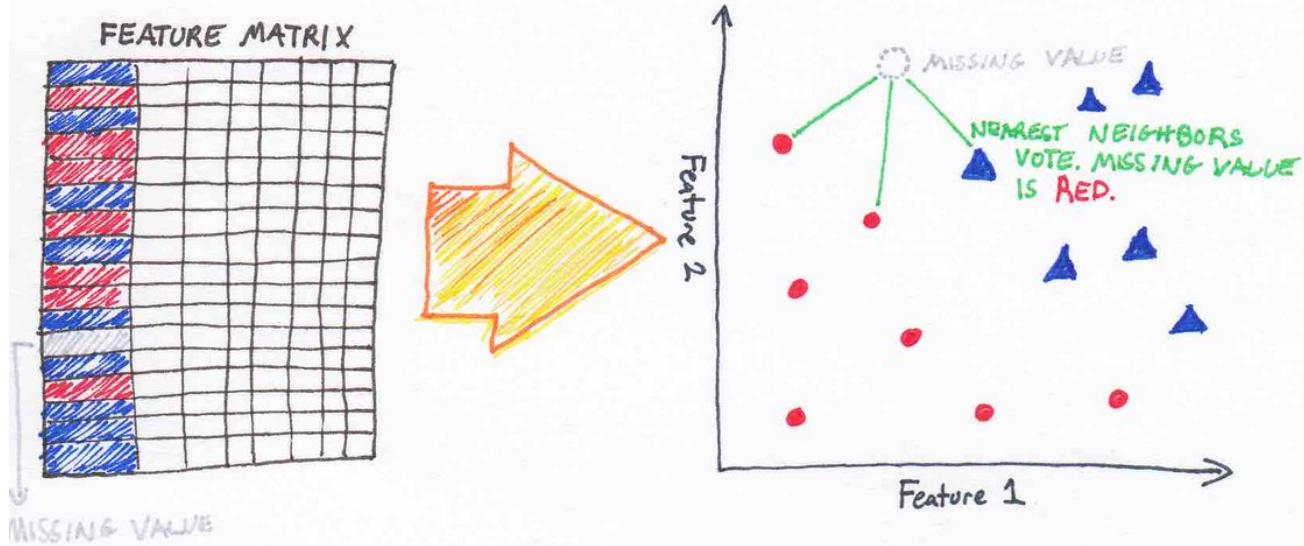
In an n -dimensional space, a hyperplane is an $n-1$ plane that divides that space.



IMPUTATION USING K-NN



Imputation Using k-NN



k-NN has value. But has to calculate distance to every point which stops being scalable in high N or/and high dimensionality.

BY CHRIS ALDON

IMPUTING MISSING VALUES

1. If quantitative, replace with an average value.
2. If qualitative, replace with most common value.
3. Use a model to predict the missing values. For example, k-nearest neighbors.

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.

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

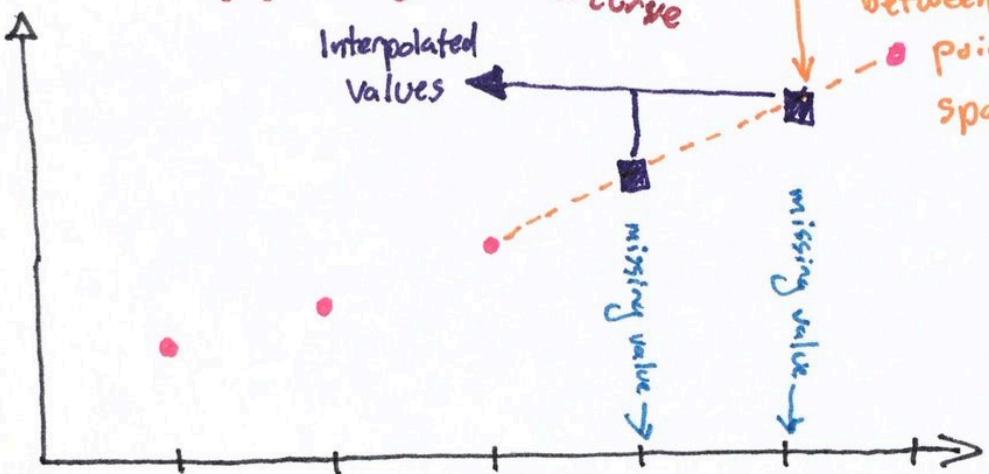
\uparrow
The interaction of features
 x_1 and x_2 .

ChrisAlbon

Interpolation

Strategy to fill in gaps of missing values by drawing a line or curve between known values

Draw a line between known points to span the gap.



JACOBIAN MATRIX

When both the inputs and outputs of a function are vectors, the matrix containing all the first-order partial derivatives is called the Jacobian.

$$\begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

Chris Albon

- Scaling is important
 - The larger the neighborhood, the less variance and more bias
 - KNN doesn't "learn" per-se. It is lazy and instead just memorizes the data.
- K-Nearest Neighbor
-
- Feature 2
- Feature 1
- class 0
- neighbors if $k=3$. Because two neighbors are green and one is red, then the unknown point is classified as green.
- class 2
- observation of unknown class
- K should be odd (so no ties)
 - Distance can be thought of as similarity
 - Variety of distance metrics can be used.
 - When we have categorical features 1 if same class and 0 otherwise. If we only have binary features we can use Hamming distance.
 - We can also weight the voting by the distance to the neighbors. For example, closer observations vote is worth more.

K-FOLD

CROSS-VALIDATION

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

OF FOLDS

K=10 IS COMMON

FOR EXAMPLE:

- MEAN SQUARED ERROR
- LOG-LOSS
- ACCURACY

K-NEAREST NEIGHBORS

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
 $K = \text{Low Bias, High Variance}$

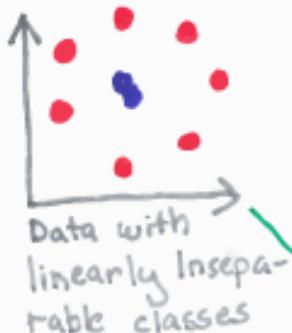


LARGE
 $K = \text{High Bias, Low Variance}$

BY CHRIS ALBON

KERNELPCA

Linear PCA will
reduce dimensionality
but not make it
linearly separable.



KPCA can reduce
dimensionality while
making data linearly
separable



Chris Albon

KERNEL TRICK

Support vector classifiers can be written as
a dot product:

$$b + \sum_{i=1}^n \alpha_i x^T x^{(i)}$$

bias → b
 α_i → parameters
 $x^T x^{(i)}$ → dot product
 $x^{(i)}$ → observation

The kernel trick is to replace the dot product with a
Kernel:

$$b + \sum \alpha_i k(x, x^{(i)})$$

$k(x, x^{(i)})$ → Kernel

Allows for non-linear decision boundaries and computational efficiency.

K-nearest Neighbors

$$\Pr(Y=j | X=x_0) = \frac{1}{K} \sum_{i \in N_0} I(y_i=j)$$

Diagram annotations:

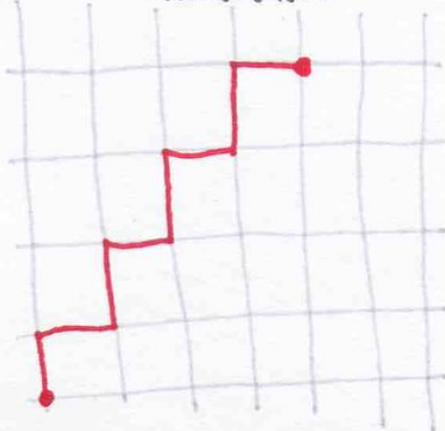
- Downward arrow from $Y=j$ to "Category j".
- Downward arrow from $X=x_0$ to "Features".
- Downward arrow from $I(y_i=j)$ to "# of neighbors".
- Downward arrow from $y_i=j$ to "True y category j".
- A green bracket labeled "indicator" covers the $I(\cdot)$ term.

L₁ norm

(Manhattan Norm)

Also Called
"Taxicab Norm"

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

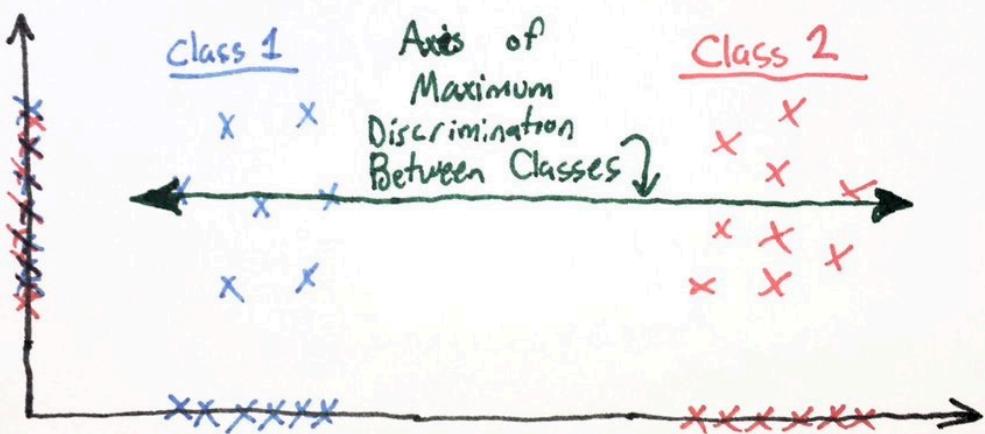


BY CHRIS ALBON

Why is nearest neighbor Lazy?

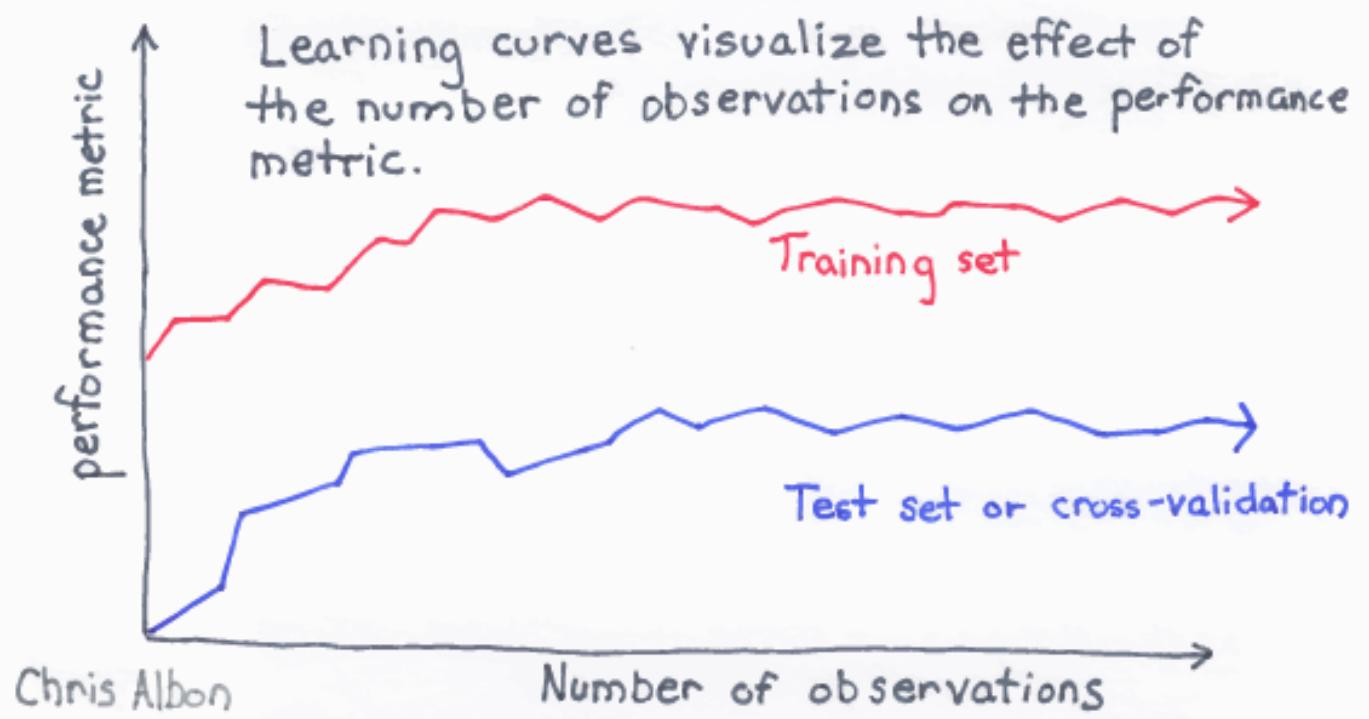
Because KNN doesn't learn so much as memorizes the data. That is, it doesn't learn any parameters.

Linear Discriminant Analysis For Dimensionality Reduction



Find Axes That Maximize Separability Between Classes.
Project data down.

LEARNING CURVE



Leave One Out Cross Validation

$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n MSE_i$$

In LOOCV, you only have sample, so it ~~is~~ is not really "Mean" squared error so much as just squared errors.

number of observation

Our score/loss metric can be anything but MSE is common in regression.

Thanks Sebastian Raschka (Buy his book)

LINEAR COMBINATION OF A SET OF VECTORS

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

Scalar

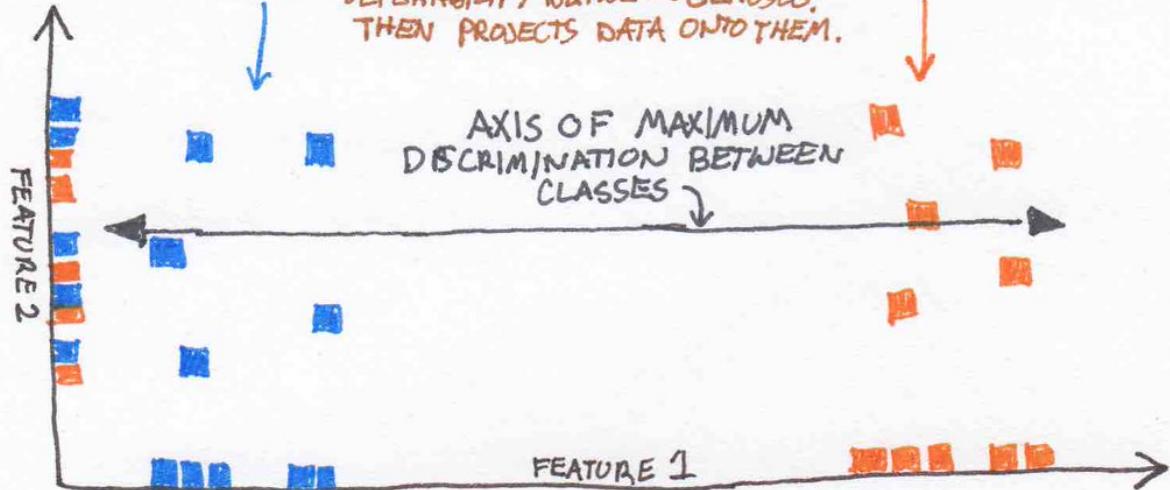
Set of vectors

ChrisAlbon



FOR DIMENSIONALITY REDUCTION

CLASS 1 FINDS AXES THAT MAXIMIZE SEPARABILITY BETWEEN CLASSES.
THEN PROJECTS DATA ONTO THEM.



Linear Regression Problems and Solutions

1. Non-linearity of X-y relationship.

DIAGNOSIS: RESIDUAL PLOT

SOLUTION: NON-LINEAR TRANSFORMATIONS OF X

2. Correlation of error terms

DIAGNOSIS: correlation

SOLUTION: EXPERIMENTAL DESIGN.

3. Non-constant error variance:

DIAGNOSIS: RESIDUAL PLOT

SOLUTION: CONCAVE TRANSFORMATION like LOG

4. Outliers:

DIAGNOSIS: RESIDUAL PLOT (STUDENTIZE) $\rightarrow \frac{e_i}{s_{e_i}}$ ← error
error's

SOLUTION: REMOVE, LOG.

5. High leverage points:

DIAGNOSIS: LEVERAGE STATISTIC

$$h_i = \frac{1}{n} + \frac{(x_i - \bar{x})}{\sum (x_i - \bar{x})^2}$$

6. Collinearity:

DIAGNOSIS: VIF

SOLUTION: DROP VARIABLE

LOG-SUM-EXP

Imagine we want to calculate:

$$c = \log \sum_{i=1}^n e^{x_i}$$

If the scale of x is very small or large it can create a number, C , that is too small or large to represent on the computer. This is called underflow and overflow. The log-sum-exp trick exploits the fact that:

$$\log \sum_{i=1}^n e^{x_i} = a + \log \sum_{i=1}^n e^{x_i - a} \quad \text{where } a = \max(x)$$

This allows us to shift the center so the greatest value of $(x_i - a)$ is zero, preventing underflow and overflow.

MANHATTAN DISTANCE

$$\sum_{i=1}^n |P_i - q_i|$$

Example:

$$Q = [3, 4, 5]$$

$$P = [0, 1, 2]$$

$$\text{MANHATTAN DISTANCE} = |0-3| + |1-4| + |2-5|$$

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

ChrisAlbon

MATRIX INVERSE

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

Inverse of matrix A Matrix n-dimensional identity matrix

ChrisAlbon

MAXNORM

The element of the vector with the largest absolute value.

$$\|x\| = \max_i |x_i|$$

vector x value i

By Chris Albon

Commonly used
in regression.

Mean Squared Error

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

true y Predicted y

$$MSE(\hat{\theta}) = \text{Var}(\hat{\theta}) + \text{Bias}(\hat{\theta}, \theta)^2$$

n
number of observations

↓
MSE as the bias-variance trade-off

MEANSHIFT CLUSTERING BY ANALOGY

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

Min-Max Scaling

$$X = \frac{X - \text{MIN}}{\text{MAX} - \text{MIN}}$$

sklearn.preprocessing. MinMaxScaler.

MIN MAX SCALING

Rescales feature values to between 0 and 1



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

Annotations for the formula:

- Original value: x_i
- Minimum value in feature: $\min(x)$
- Maximum value in feature: $\max(x)$

Chris Albon

MAR

MISSING AT RANDOM

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

For example, men are less likely to answer a salary question in a survey but we capture gender identity in another question.

Chris Albon

MNAR

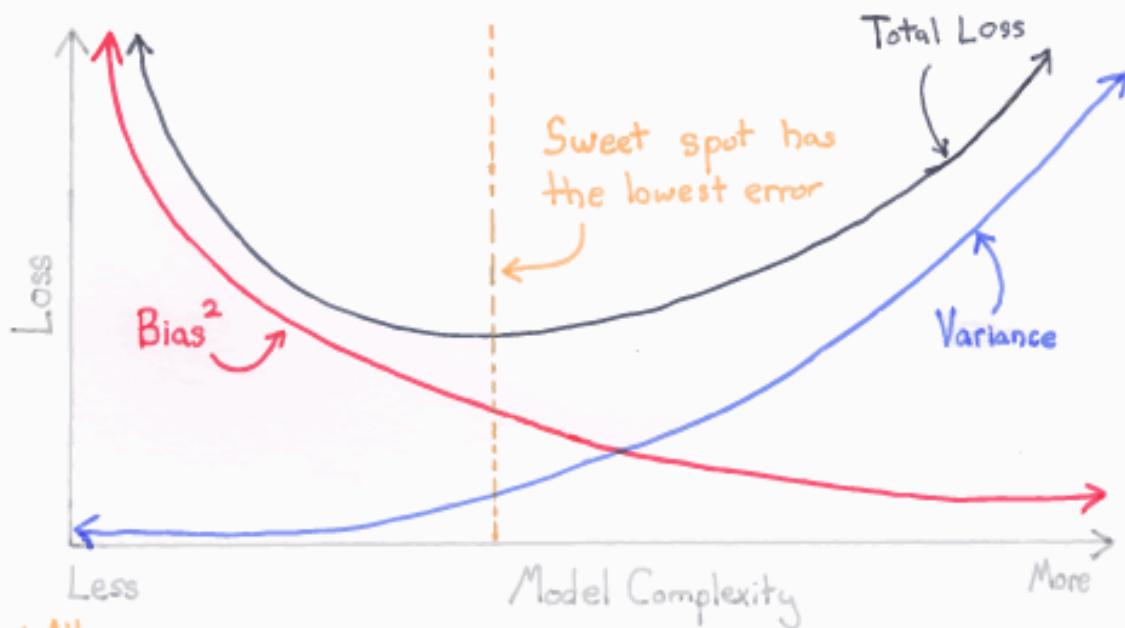
MISSING NOT AT RANDOM

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.

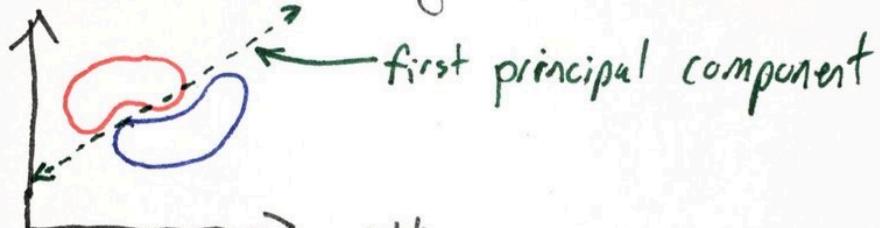
Chris Albon

MODEL COMPLEXITY



Motivation For Kernel PCA

When data is not linearly separable:

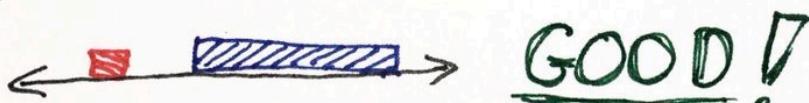


PCA (linear) will ^{might} merge classes:



BAD!

But Kernel PCA can work:



GOOD!

MSE vs MAE

Benefits of Squaring:

- Always returns a positive value
- Emphasizes large errors ← important
- Continuously differentiable
- Version of L₂ norm

BY CHRIS ALBON

Multi colinearity

When collinearity exists between three or more coefficients

Test with VARIANCE INFLATION FACTOR (VIF)

$$\frac{\text{Var}(\hat{B}_j)}{\text{Var}(\hat{B}_j)} \text{ when part of the whole model}$$

VIF of $\underline{5}$ or ~~more~~ is considered bad.

NAIVE BAYES CLASSIFIER!

$$p(\text{class} \mid \text{data}) = \frac{p(\text{data} \mid \text{class}) \times p(\text{class})}{p(\text{data})}$$

Normal distribution

$$\frac{1}{\sqrt{2\pi \text{Var}(\text{height})}} e^{-\frac{(x_i - \bar{x})^2}{2\pi \text{Var}(\text{height})}}$$

$$\text{posterior (male)} = \frac{p(\text{height} \mid \text{male}) p(\text{weight} \mid \text{male}) \times p(\text{male})}{\text{marginal probability}}$$

We ignore this, only compare numerators of the posterior!

BY CHAIS ALBON

NATURAL LOG

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

Euler's number: 2.71828...

BY CHRIS ALBON

NON-PARAMETRIC METHODS

Non-parametric methods do not assume a functional form of the relationship between X and y .

ADVANTAGE: Can fit a wider range of types of relationships.

DISADVANTAGE: Often a large number of observations required to fit a model of any quality.

BY CHRIS ALBON

Normalization

Rescales the values of individual observation so that they ~~are~~ have unit (1) norm.

Note: A lot of ambiguity around "normalization". Use scikit's lexicon.

Option 1: Euclidean Norm

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

Option 2: Manhattan Norm

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

Seen in text classification

NOTATION 1

¬

Not

∃

Such that

$a \in b$ a is an element
of b ∴ Therefore

Chris Albon

NATIONS OF PROBABILITY

Frequentist: The rate at which an event will occur if we repeat the experiment.

Bayesian: The degree of belief that something is true.

Chris Albon

Odds and Log Odds
in Logistic Regression

$$\text{ODDS: } \frac{p(x)}{1-p(x)} = e^{B_0 + B_1 x_1 + \dots + B_n x_n}$$

$$\text{LOG ODDS: } \log\left(\frac{p(x)}{1-p(x)}\right) = B_0 + B_1 x_1 + \dots + B_n x_n$$

ODDS

$$\frac{\Pr(y)}{\Pr(\sim y)}$$

Odds is the ratio of the probability an event occurs with the probability of an event not occurring.

ChrisAlbon

One-Hot Encoding

Feature
Apple
Pear
Apple
Pear
Apple



	Apple	Pear
Apple	1	0
Pear	0	1
Apple	1	0
Pear	0	1
Apple	1	0

Note:
Splitting one feature into many binary features can artificially reduce the feature importance in random forest because info is spread over many features.

One-Hot encoding allows us to turn categorical data into features with numerical values (that many ML algorithms need) while not mathematically imply any ordinal relationship between the classes.

OUT-OF-BAG ERROR

Method of estimating test error in bagged models. OOB error is the mean error of all observations, x_i , using only the trees which did not use x_i in their bootstrapped datasets. OOB error is as accurate as a test set of equal size as the training set (Breiman 1996).

BY CHRIS ALBON

OUTLIARS

DROP : Not a great option. We lose lots of info.
Find out if genuine extreme value or broken sensor.

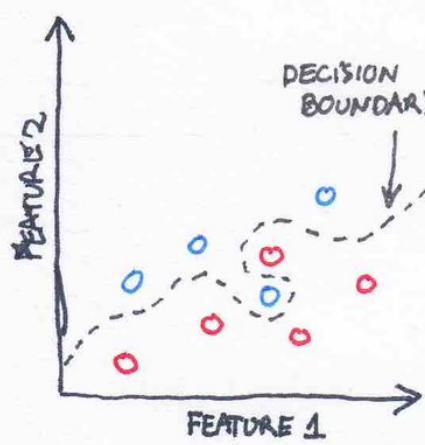
MARK

- Safest option. We can see if the outliers had an effect.

RESCALE

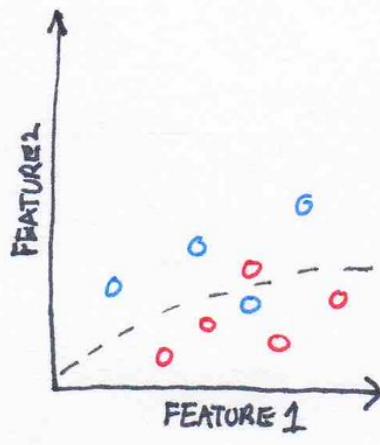
- Log values so outliers don't have as great an effect.

OVERFIT vs UNDERFIT

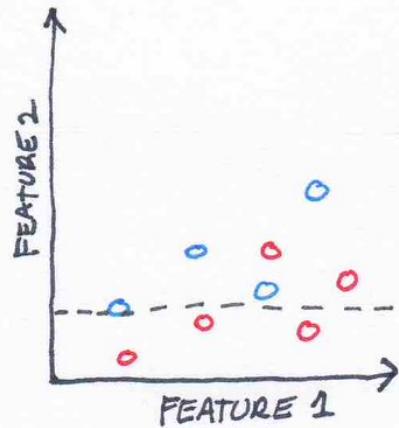


OVERFIT
"HIGH VARIANCE"

BY CHRIS ALBON



IDEAL



UNDERFIT
"HIGH BIAS"

Parameters vs. Hyperparameters

PARAMETERS: LEARNED FROM THE TRAINING. FOR EXAMPLE: REGRESSION COEFFICIENTS.

HYPERPARAMETERS: SET BEFORE THE TRAINING STEP. "TUNED" THROUGH HYPERPARAMETER GRID SEARCH AND RELATED METHODS.

Parametric Methods

- ① Assume a functional form of f . For example: we can assume that f is linear.
- ② Use a method (OLS, MLE, etc) to train the model.

ADVANTAGE: If we choose the correct functional form, we ~~need~~ the estimation ~~problem~~ problem to ^{reduce} estimating a small set of parameters.

Pearson's Correlation Coefficient

$$P = \frac{\text{Cov}(X, Y)}{S_x S_y}$$

← covariance
← standard deviations

GOOD IF: Variables are roughly normal or if the relation of the variables is roughly linear.

BAD IF: Outliers present or skewed distribution.

CORRELATION

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

(Pearson's R) 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.

BY CHRIS ALGON

PERCEPTRON LEARNING

$$\hat{w}_j \equiv w_j + \Delta w_j$$

Updated parameter
Current parameter
Learning rate
 $\Delta w_j = \eta (y_i - \hat{y}_i) x_{ij}$

True y_i Input value
Predicted \hat{y}_i

ChrisAlbon

Polynomial Regression

$$\hat{y}_i = B_0 + B_1 x_i + B_2 x_i^2 + \dots + e_i$$

ADDS NON-LINEARITY
WHILE STILL OBEDIING THE
RULES OF A LINEAR FUNCTION.

POWER RULE

What is $\frac{d}{dx} x^3$? Use the power rule:

$$\frac{d}{dx} x^n = nx^{n-1}$$

Therefore the derivative is $3x^2$

ChrisAlban

Precision - Recall Trade-off

Low Precision

High Recall

High Precision

Low Recall

Optimistic Model

(Predict almost
everything as true)

Pessimistic
Model

(only predict positives
which very true)

Use the probability threshold
to move between pessimism
and optimism.

Precision

"Precision is about the predicted positives"

$$\text{Precision} = \frac{\text{TruePositives}}{\text{TruePositives} + \text{FalsePositives}}$$

$$\text{Precision} = \frac{\text{True Positives}}{\text{Everything Predicted Positive}}$$

Precision is the ability of a classifier not to label as positive an observation that is negative. It measures the purity of positive predictions.

BY CHRIS ALBON

PREPROCESSING TEST + TRAINING SETS

1. Fit preprocessor to the training set ONLY.
2. Apply it to both the training and test set.

WHY?

Because we have to "pretend" the test set is unknown data.

BY CHRIS ALBON

Preprocessing When Tests and Training.

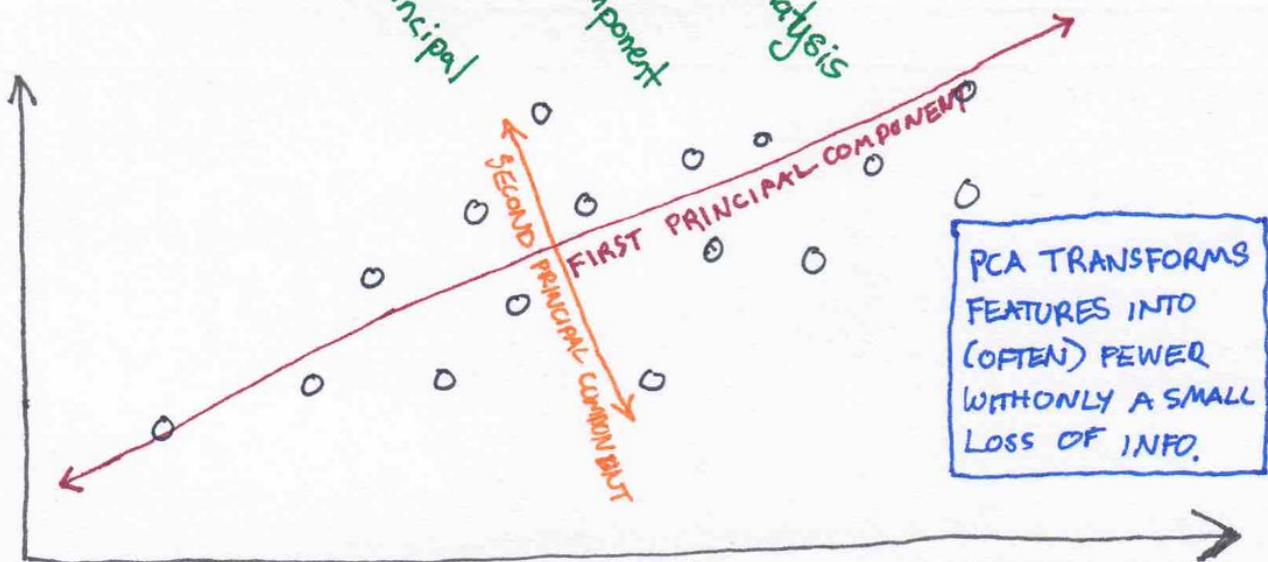
- 1) Fit the preprocessor to the training set.
- 2) Apply it to both the training AND test.set.

Why? Because we have to pretend the test set is new data.

PCA

Principal Component Analysis

PCA PROJECT DATA
ONTO THE PRINCIPAL'S
COMPONENTS.



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.

ChrisAlbon

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



R²

R² looks at
How much variance
in the target
vector is explained
by the features

$$R^2 = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$

TRUE Y PREDICTED Y

VARIANCE IN PREDICTIONS VS. TRUE Y_i

VARIANCE IN TARGET VECTOR

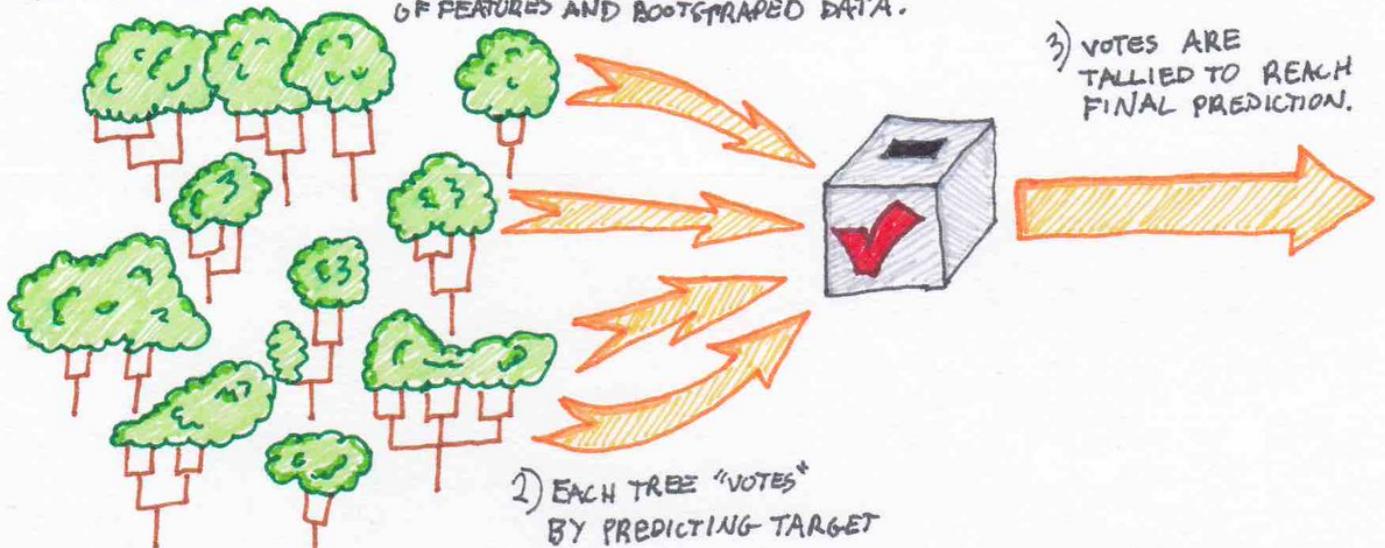
TRUE Y MEAN TRUE Y

BY CHRIS ALBON

RANDOM FOREST

1) MANY TREES ARE CREATED USING RANDOM SUBSET OF FEATURES AND BOOTSTRAPPED DATA.

CLASSIFICATION



BY CHRIS ALBON

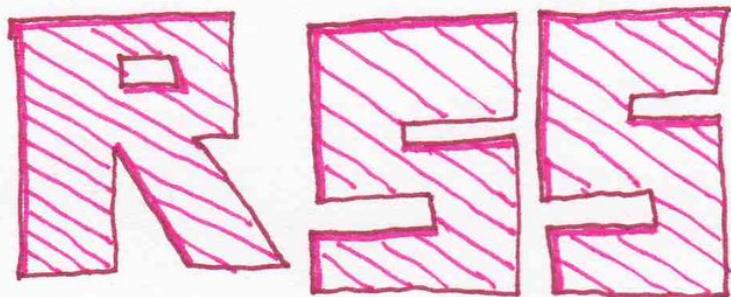
recall

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

"Recall is about the real positives"

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

BY CHRIS ALBON



Residual

Sum of Squares

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

true y_i

\hat{y}_i → predicted y_i
(Also called "y-hat")

Note that squaring the error with more heavily penalty a few large errors over many small errors even if the sum of errors are the same

RIDGE REGRESSION

Residual sum of squares →

RSS

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

Tuning parameter

Shrinkage

Parameters squared

Remember:
Standardize
the data first.

Chris Albon

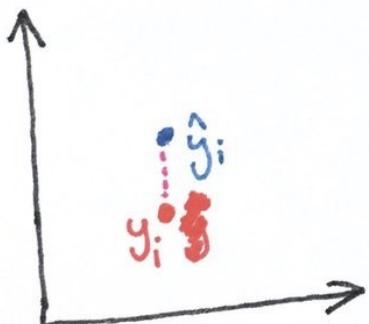
Disadvantage:
parameters cannot
be zero like
with Lasso
regression.

RSS
Residual Sum Of Squares

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

true y value \hat{y}_i = predicted y

Removes negative signs,
 penalizes larger errors
 more.



example!

y_i	\hat{y}_i	$(y_i - \hat{y}_i)^2$	RSS
2	2	$0^2 = 0$	
4	2	$2^2 = 4$	
-4	-2	$-2^2 = 4$	
3	1	$2^2 = 4$	8

SATURATION

When a function's output is very insensitive to inputs. Example: Sigmoid function.



Selecting # of components in PCA

Select smallest number of components such that:

$$\frac{\frac{1}{m} \sum_{i=1}^m \|x_i - x_i \text{ appox}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x_i\|^2} = 0.01$$

DISTANCE BETWEEN X AND ITS PROJECTIONS. ← AVERAGE SQUARE PROJECTION ERROR
99% of VARIANCE RETAINED ← TOTAL VARIATION IN DATA

Sensitivity

Recall on real positive examples.

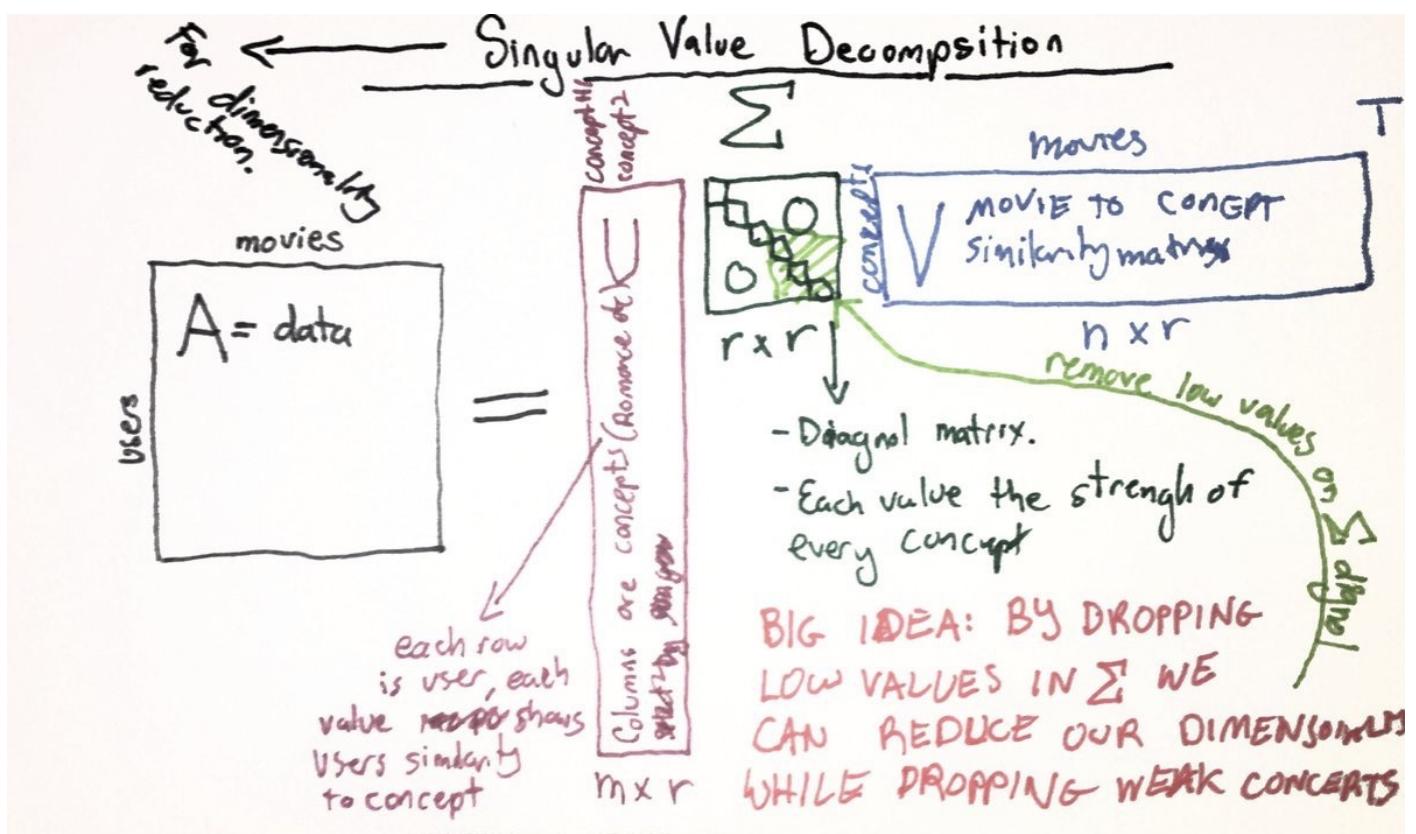
$$\frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

Everything predicted positive correctly {

True Positives

{ True Positives + False Negatives

Everything actually positive

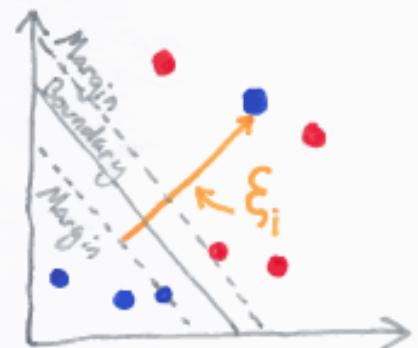


SLACK VARIABLE IN SOFT-MARGIN SVM

$\xi_i = 0$, the observation is correctly classified

$0 < \xi_i < 1$, the observation is inside the margins

$\xi_i \geq 1$, the observation is misclassified



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

SPARSITY

Sparse matrices allow for more efficient data storage, particularly when most values are zero or missing.

$$\text{Sparsity} = \frac{\text{\# of zero/missing values}}{\text{Total \# of values}}$$

$$\text{Sparsity} = 1 - \text{density}$$

BY CHAIS ALBON

SQUARE ROOT

Principle Root

$$\sqrt[r]{x} = x^{\frac{1}{r}}$$

Square Root

$$\sqrt[2]{x} = x^{\frac{1}{2}}$$

Often used to reverse
or undo a Squaring.
For example:

$$\sqrt{x^2} = x$$

BY CHRIS ALBON

Standard Error

ANSWERS: TELLS US THE AVERAGE AMOUNT $\hat{\mu}$ DIFFERS FROM THE TRUE μ .

$$\text{Variance} = \frac{\sum_i (x - \mu)^2}{N}$$

$$\text{Var}(\hat{\mu}) = \text{SE}(\hat{\mu})^2 = \frac{\sigma^2}{n}$$

\downarrow \downarrow \downarrow
Sample mean Sample mean # of observations

unit is ~~#~~ the standard deviation

Standardization

(z-score scaling)

Rescales input data (features) to be approximately standard normally distributed.

$$\text{Scaled value } \leftarrow \bar{x}' = \frac{\text{value } x - \bar{x}}{\sigma}$$

mean value
of feature

\hookrightarrow standard deviation

STOCHASTIC GRADIENT DESCENT

1. Shuffle Observations.

2. For each observation:

$$\hat{w} = w - \eta \nabla \text{Loss}(w, x_i, y_i)$$

Diagram illustrating the update rule:

- \hat{w} : updated parameter
- w : parameter
- η : Learning Rate
- $\nabla \text{Loss}(w, x_i, y_i)$: Gradient
- x_i, y_i : Individual x and y value
- w : Parameter

3. Repeat until ~minimum value achieved.

By CHRIS ALBON

STOP WORDS

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

Examples

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

STRATEGIES WHEN YOU HAVE **HIGH VARIANCE**

- Weight Decay
- Drop out or bagging
- Dimensionality reduction
- Feature selection

Chris Albon

SUPERVISED VS. UNSUPERVISED

In supervised learning, for every observation we have some feature values x_i and some target vector, y_i . We use both to train a model that take in some x_i values and outputs a predicted \hat{y}_i .

In unsupervised learning, we only have ~~the~~ features ~~values~~, and do not have a target vector. This makes the estimation problem much more difficult. When possible, use supervised learning.

BY CHRIS ALBON

SUPERVISED VS UNSUPERVISED

In supervised learning, for every observation we have some feature values and some target vector or tensor. We use both to train a model that takes in some x values and outputs a predicted value.

In unsupervised learning, we only have features and do not have a target. This makes the estimation problem more difficult. When possible, use supervised learning.

Chris Albon

T-STATISTIC

Number of standard deviations a parameter is away from a constant.

$$t = \frac{\hat{\beta}_i - C}{SE(\hat{\beta}_i)}$$

Parameter
Some constant.
Often zero.

Standard error

Chris Albon

T - Tests

One-sample t-test: Compare sample to some constant.

$t \text{ statistics} \leftarrow$ mean $\leftarrow \bar{x} - \mu_0 \leftarrow \text{some value}$
standard deviation $\leftarrow s / \sqrt{n} \leftarrow \# \text{ of obs}$

Unpaired t-test: Compare if two distributions are different



Paired t-test: Compare if two distributions of repeated sampling are the same.

TF-IDF

THIS IS THE
NOTATION USED
BY SCIKIT-LEARN

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

number of times term appears in a doc. \uparrow

of documents where a word appears \uparrow

term frequency \downarrow

of times word appear in a doc

inverse document frequency \downarrow

$\log \frac{1 + \frac{n_d}{\text{df}(d, f)}}{1 + \frac{1}{n_d}} + 1$

number of documents \downarrow

document frequency of the term, t . \downarrow

BY CHRIS ALBON

The Argument For Parametric Models

- When the data generating function roughly matches a parametric probability distribution we can limit our calculations to only its parameters (mean, variance, etc.). This lets us know (assume) a lot using only a little information.
- The flexibility of many probability distributions means that often the ~~choice~~ choice of the functional family of distribution (normal vs. student etc.) is ~~only~~ ^{often} not a problem
- Range matching is, however, important. If we want a probability, we should not choose a probability distribution that outputs numbers greater than 1, for example.

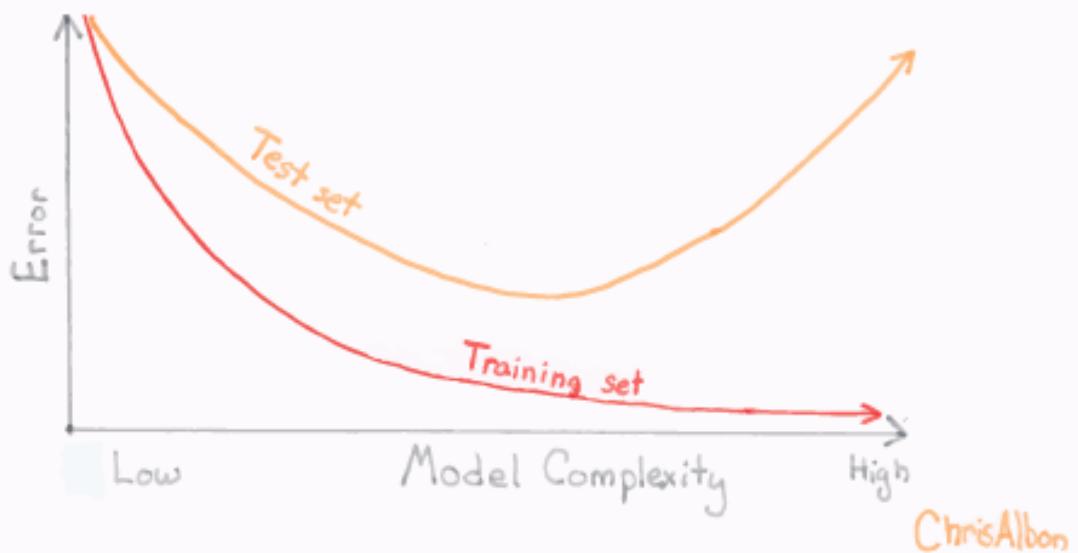
BY CHRIS ALBON

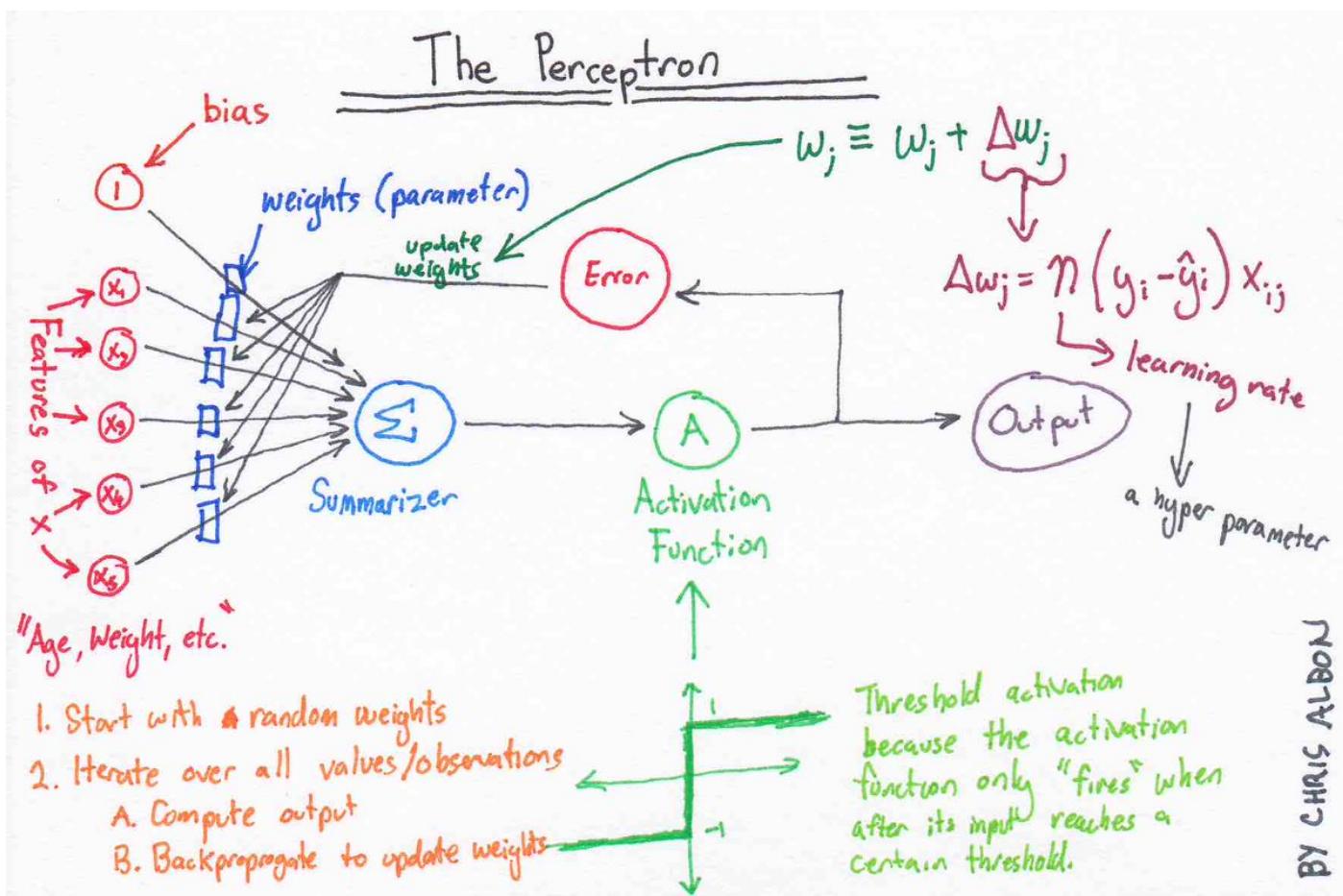
THE EFFECT OF **DROPOUT** ON HIDDEN UNITS

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





THE RANDOM IN RANDOM FOREST

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

BY CHRIS ALBON

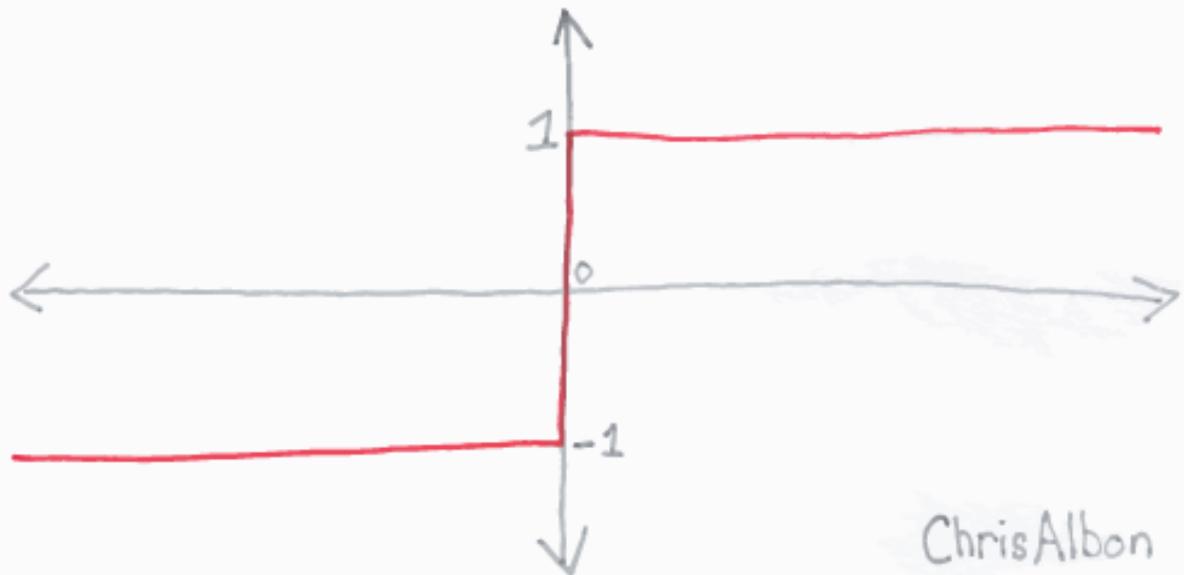
THE THEREFORE
BECAUSE

NOTATION

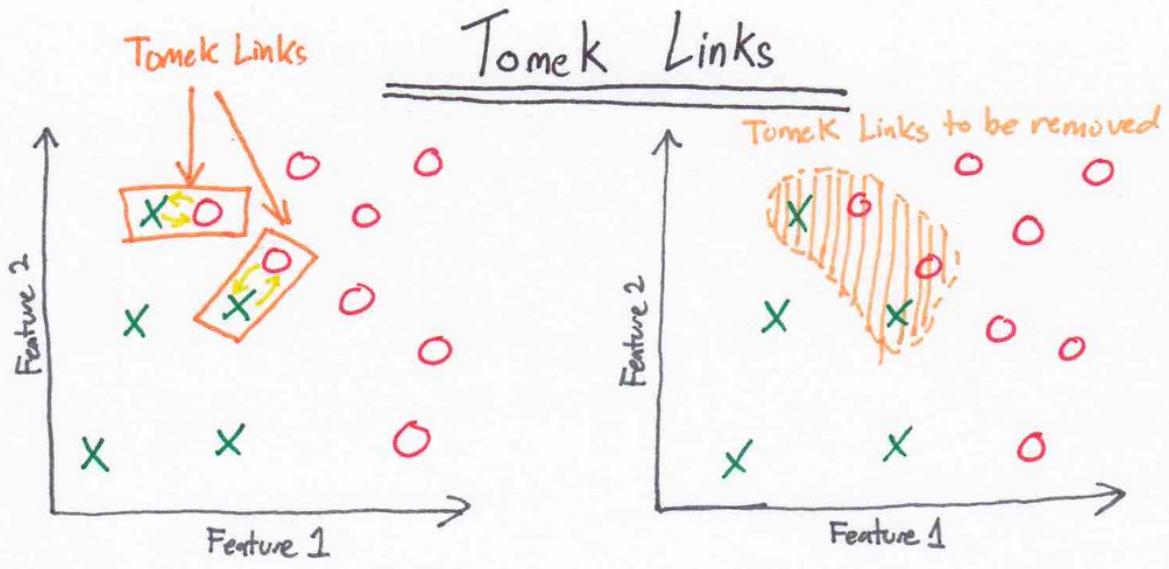
+ THEREFORE
because

BY CHRIS ALBON

THRESHOLD ACTIVATION



Chris Albon



Tomek Link:

1. x 's nearest neighbor is y
2. y 's nearest neighbor is x
3. x and y are different classes.

Removing Tomek links has been shown to improve model performance. But problematic when n is small or classes are highly imbalanced.

BY CHRIS ALBON

Training Error Rate

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

True y_i

Predicted \hat{y}_i

Indicator

Failed Prediction

All Predictions

number of observations

Training ERROR RATE

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

↓
number of observations

↓ indicator

↓ True y_i

↓ Predicted \hat{y}_i

FAILED PREDICTIONS
ALL PREDICTIONS

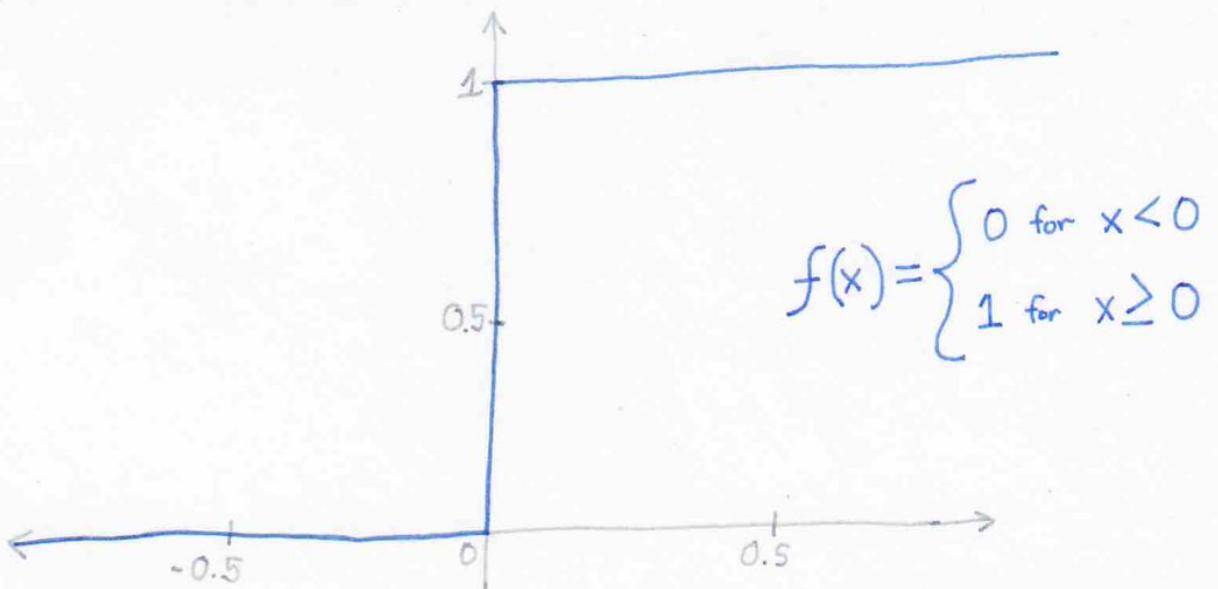
BY CHRIS ALBON

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

UNIT-STEP ACTIVATION FUNCTION



BY CHRIS ALBON

VANISHING GRADIENT PROBLEM

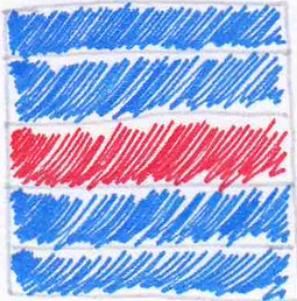
When the gradients of the loss function with respect to the parameters in the early layers of a network are very small. Causes slow learning and since many gradients are tiny, they don't contribute much to learning and can lead to poor performance

Chris Albon

VARIANCE FOR FEATURE SELECTION

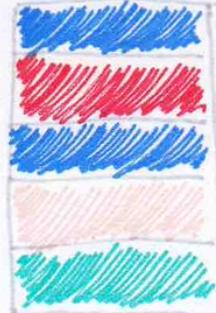
VARIANCE IS A USEFUL PROXY FOR THE INFORMATION CONTAINED IN A FEATURE. MORE VARIANCE MEANS MORE INFORMATION. AND THIS MORE USED IN TRAINING.

FEATURE 1, LOW VARIANCE



WITH LOW VARIANCE,
THIS FEATURE HAS LESS
ABILITY TO TRAIN A
MODEL.

FEATURE 2, HIGH VARIANCE



WITH HIGH VARIANCE
THIS FEATURE HAS
MORE INFO TO
TRAIN A MODEL.

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

ChrisAlbon

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.

BY CHRIS ALBON

What are principal components?

PCA finds low dimensional representation of data that contains as much of the variance of the original data as possible.

↳ a measure of interestingness

Principal components are the linear combination of features that has the maximum variance out of all linear combinations that are uncorrected to the ~~and~~ previous principal component.

Alternative interpretation: principal components ~~possible~~ are low dimensional linear surfaces closest to the observations.

Youden's J statistic

$$J = \frac{\text{How good are we at predicting positives}}{TP + FN} + \frac{\text{How good are we at predicting negatives.}}{TN + FP} - 1$$

Values range between -1 and 1 . 0 = test is useless (same number of false positives as false negatives). 1 = test is perfect (no false positives and no false negatives). Often used to select the cut-off in a ROC. Highest J is best cut-off.

Z-score

$$\frac{\bar{x} - \mu_{\bar{x}}}{\sigma_{\bar{x}}}$$

I'm drunk!

Z-score is the number of standard deviations away from the mean.

ZERO-ONE LOSS

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

Indicator Function

True class y_i Predicted class \hat{y}_i True class y_i

BY CHRIS ALBON